

Transparent numerical boundary conditions for evolution equations: Derivation and stability analysis

Jean-François COULOMBEL*

September 23, 2016

Abstract

The aim of this article is to propose a systematic study of transparent boundary conditions for finite difference approximations of evolution equations. We try to keep the discussion at the highest level of generality in order to apply the theory to the broadest class of problems.

We deal with two main issues. We first derive transparent numerical boundary conditions, that is, we exhibit the relations satisfied by the solution to the pure Cauchy problem when the initial condition vanishes outside of some domain. Our derivation encompasses discretized transport, diffusion and dispersive equations with arbitrarily wide stencils. The second issue is to prove sharp *stability* estimates for the initial boundary value problem obtained by enforcing the boundary conditions derived in the first step. We focus here on discretized transport equations. Under the assumption that the numerical boundary is non-characteristic, our main result characterizes the class of numerical schemes for which the corresponding transparent boundary conditions satisfy the so-called Uniform Kreiss-Lopatinskii Condition introduced in [GKS72]. Adapting some previous works to the non-local boundary conditions considered here, our analysis culminates in the derivation of trace and semigroup estimates for such transparent numerical boundary conditions. Several examples and possible extensions are given.

AMS classification: 65M06, 65M12, 35L02, 35K05, 35Q41.

Keywords: evolution equations, difference approximations, transparent boundary conditions, stability.

Throughout this article, we use the notation

$$\mathcal{U} := \{\zeta \in \mathbb{C}, |\zeta| > 1\}, \quad \mathbb{D} := \{\zeta \in \mathbb{C}, |\zeta| < 1\}, \quad \mathbb{S}^1 := \{\zeta \in \mathbb{C}, |\zeta| = 1\}, \\ \overline{\mathcal{U}} := \mathcal{U} \cup \mathbb{S}^1, \quad \overline{\mathbb{D}} := \mathbb{D} \cup \mathbb{S}^1.$$

We let $\mathcal{M}_{n_1, n_2}(\mathbb{K})$ denote the set of $n_1 \times n_2$ matrices with entries in $\mathbb{K} = \mathbb{R}$ or \mathbb{C} . In the case $n_1 = n_2 = n$, we use the notation $\mathcal{M}_n(\mathbb{K})$ for the set of square matrices of size n . If $M \in \mathcal{M}_n(\mathbb{C})$, M^* denotes the conjugate transpose of M and $\text{sp}(M)$ denotes the spectrum of M . We let I denote the identity matrix or the identity operator when it acts on an infinite dimensional space. The subscript in I_k is intended to make the dimension k of the underlying vector space \mathbb{C}^k precise when needed. We use the notation

*CNRS and Université de Nantes, Laboratoire de Mathématiques Jean Leray (UMR CNRS 6629), 2 rue de la Houssinière, BP 92208, 44322 Nantes Cedex 3, France. Email: jean-francois.coulombel@univ-nantes.fr. Research of the author was supported by ANR project BoND, ANR-13-BS01-0009-01.

$x^* y$ for the Hermitian product $\sum_i \bar{x}_i y_i$ of two vectors $x, y \in \mathbb{C}^n$. For two vectors $x, y \in \mathbb{C}^n$, the quantity $\sum_i x_i y_i$ is denoted $x \cdot y$; it coincides with the Euclidean product when the vectors have real coordinates. The norm of a vector $x \in \mathbb{C}^n$ is $|x| := (x^* x)^{1/2}$. The induced matrix norm on $\mathcal{M}_n(\mathbb{C})$ is denoted $\|\cdot\|$.

The letter C denotes a constant that may vary from line to line or within the same line. The dependence of the constants on the various parameters is made precise throughout the text. If a constant C depends on some parameter v , we write either C_v or $C(v)$ to make this dependence explicit.

In what follows, we let $d \geq 1$ denote a fixed integer, which will stand for the dimension of the space domain \mathbb{R}^d or \mathbb{Z}^d we are considering. We shall use the space ℓ^2 of square integrable sequences. Sequences may be valued in \mathbb{C}^k for some integer k . In that case, we write $\ell^2(\mathbb{Z}^d; \mathbb{C}^k)$ to emphasize that sequences are vector valued. Some sequences will be indexed by \mathbb{Z}^{d-1} while some will be indexed by \mathbb{Z}^d or a subset of \mathbb{Z}^d . We thus introduce some specific notation for the norms. Let $\Delta x_i > 0$ for $i = 1, \dots, d$ be d space steps. We shall make use of the $\ell^2(\mathbb{Z}^{d-1})$ norm that we define as follows: for all $v \in \ell^2(\mathbb{Z}^{d-1})$,

$$\|v\|_{\ell^2(\mathbb{Z}^{d-1})}^2 := \left(\prod_{k=2}^d \Delta x_k \right) \sum_{i=2}^d \sum_{j_i \in \mathbb{Z}} |v_{(j_2, \dots, j_d)}|^2.$$

The corresponding scalar product is denoted $\langle \cdot, \cdot \rangle_{\ell^2(\mathbb{Z}^{d-1})}$. Then for all integers $m_1 \leq m_2$, we set

$$\|u\|_{m_1, m_2}^2 := \Delta x_1 \sum_{j_1=m_1}^{m_2} \|u_{(j_1, \cdot)}\|_{\ell^2(\mathbb{Z}^{d-1})}^2,$$

to denote the ℓ^2 norm on the set $[m_1, m_2] \times \mathbb{Z}^{d-1}$ (m_1 may equal $-\infty$ and m_2 may equal $+\infty$). The corresponding scalar product is denoted $\langle \cdot, \cdot \rangle_{m_1, m_2}$. In the particular case $d = 1$, the space step is denoted Δx and the ℓ^2 norm on the interval $[m_1, m_2]$ reduces to

$$\|u\|_{m_1, m_2}^2 := \Delta x \sum_{j=m_1}^{m_2} |u_j|^2.$$

The $\ell^2(\mathbb{Z}^{d-1})$ norm reduces to the norm of vectors. Other notation is introduced when needed throughout the text or is meant to be self-explanatory.

1 Introduction

1.1 The context

We are concerned here with the approximation of partial differential equations of the evolutionary type in the whole space \mathbb{R}^d . For simplicity, we restrict here to linear partial differential equations with constant coefficients of the form

$$\partial_t v + P(\partial_1, \dots, \partial_d) v = 0, \tag{1}$$

where the differential operator P is a polynomial expression of the spatial partial derivatives $\partial_1, \dots, \partial_d$ (we write ∂_j for the partial derivative with respect to the j -th space variable x_j and use ∂_t for the partial derivative with respect to time). The differential operator P may have complex coefficients so that the above framework encompasses the Schrödinger equation, as well as prototype evolution equations for real valued functions such as the transport, heat or Airy equation. We restrict here for simplicity to the case of *scalar* evolution equations of order 1 in time: the unknown v in (1) is either real or complex valued

and there is no second or higher order time derivative in (1). There would be no great effort to consider higher order equations such as the wave or beam equations but the functional framework would be slightly different.

When considered on the whole space domain \mathbb{R}^d , solving (1) usually relies on the Fourier transform and we assume that well-posedness holds for (1) in $L^2(\mathbb{R}^d)$. More precisely, we assume

$$\forall \boldsymbol{\xi} \in \mathbb{R}^d, \quad \operatorname{Re} P(i \xi_1, \dots, i \xi_d) \geq 0,$$

so the solution to the Cauchy problem (1) satisfying $v|_{t=0} = v_0$ reads

$$\forall t \geq 0, \quad v(t, x) = \frac{1}{(2\pi)^d} \int_{\mathbb{R}^d} e^{ix \cdot \boldsymbol{\xi}} e^{-t P(i \xi_1, \dots, i \xi_d)} \widehat{v}_0(\boldsymbol{\xi}) d\boldsymbol{\xi}, \quad \widehat{v}_0(\boldsymbol{\xi}) := \int_{\mathbb{R}^d} e^{-ix \cdot \boldsymbol{\xi}} v(x) dx.$$

Of course, for hyperbolic or dispersive equations, $P(i \xi_1, \dots, i \xi_d)$ is a purely imaginary number and the solution v to (1) is also defined for $t \leq 0$. However, the theory developed below for numerical schemes is mostly restricted to evolution equations in positive times because, even though the original partial differential equation (1) may be time reversible, most of its finite difference approximations will not be so. Hence we consider $t \geq 0$ in what follows, which is of course no restriction when dealing with parabolic equations. The above Fourier representation of the solution implies by Plancherel's Theorem the uniform bound

$$\forall t \geq 0, \quad \|v(t, \cdot)\|_{L^2(\mathbb{R}^d)} \leq \|v_0\|_{L^2(\mathbb{R}^d)},$$

and we shall be interested below in numerical approximations of (1) for which the same L^2 decay (or conservation) property holds, or a slightly relaxed version of it.

We now turn to the finite difference approximation of (1). We introduce a time step $\Delta t > 0$ and some space steps $\Delta x_i > 0$, $i = 1, \dots, d$. The solution v to (1) is ‘approximated’ by a piecewise constant function

$$u(t, x) := u_j^n, \quad \forall (t, x) \in [n \Delta t, (n+1) \Delta t) \times \prod_{k=1}^d [j_k \Delta x_k, (j_k + 1) \Delta x_k), \quad (2)$$

where the sequence $(u_j^n)_{n \in \mathbb{N}, j \in \mathbb{Z}^d}$ is defined as the solution to the recurrence relation:

$$\begin{cases} \sum_{\sigma=0}^{s+1} Q_\sigma u^{n+\sigma} = 0, & n \geq 0, \\ (u^0, \dots, u^s) = (f^0, \dots, f^s) \in \ell^2(\mathbb{Z}^d)^{s+1}. \end{cases} \quad (3)$$

In (3), each (real or complex) sequence u^n, \dots, u^{n+s+1} is defined on \mathbb{Z}^d , that is, we consider a pure Cauchy problem on the whole space, and the operators Q_σ are given by:

$$\forall \sigma = 0, \dots, s+1, \quad Q_\sigma := \sum_{\ell=-r}^p a_{\ell, \sigma}(\Delta t, \Delta x) \mathbf{S}^\ell, \quad \text{with } (\mathbf{S}^\ell u)_j := u_{j+\ell}. \quad (4)$$

Let us comment a little on (3), (4). First of all, the *stencil* of the scheme (3) is assumed to be *fixed* and finite. More precisely, we consider some *fixed* numbers $r_1, \dots, r_d, p_1, \dots, p_d \in \mathbb{N}$ and we use in (4) the short notation

$$\sum_{\ell=-r}^p := \sum_{k=1}^d \sum_{\ell_k=-r_k}^{p_k}.$$

Then in (4), the coefficients $a_{\ell,\sigma}(\Delta t, \Delta x)$ are either real or complex numbers (depending on whether (1) has real or complex coefficients/solutions) and they may depend on the time and space steps, which we abbreviate as $(\Delta t, \Delta x)$. Each operator Q_σ in (3) thus acts on sequences indexed by the (discrete) spatial variable $j \in \mathbb{Z}^d$, and is bounded on ℓ^2 for any given choice of the discretization parameters $(\Delta t, \Delta x)$ under consideration. The (discrete) time variable n enters as a parameter when we apply each operator Q_σ . For convenience, we write most of the time $Q_\sigma u_j^{n+\sigma}$ rather than $(Q_\sigma u^{n+\sigma})_j$ to denote the application of the operator Q_σ to the sequence $u^{n+\sigma}$, the resulting sequence being evaluated at the space index j . For simplicity, we shall not write that the Q_σ 's depend on the time and space steps $(\Delta t, \Delta x)$ even though they will in all applications we have in mind.

Our aim is to cover several situations within the same framework. In particular, we wish to cover *explicit* schemes (in that case, Q_{s+1} is the identity) that are commonly used for discretizing transport equations and which come unavoidably with CFL type restrictions [Str62], and *implicit* discretizations of diffusion and/or dispersive equations where the goal is precisely to avoid stringent stability constraints. We therefore consider from now on that the time and space steps belong to some set

$$\Delta \subset (0, +\infty)^{d+1},$$

where for instance we wish Δ to be the semi-open square $(0, 1] \times (0, 1]$ when discretizing the one-dimensional heat or Schrödinger equation by an implicit scheme, and Δ can be a semi-open interval

$$\{(\Delta t, \Delta t/\lambda_1, \dots, \Delta t/\lambda_d), \quad \Delta t \in (0, 1]\},$$

for some *fixed* parameters $\lambda_1, \dots, \lambda_d > 0$ when we deal with an explicit scheme for a transport equation in \mathbb{R}^d . In all what follows, we assume that the coefficients $a_{\ell,\sigma}$ in (4) are defined on Δ . For convergence purposes, it is tacitly assumed that Δ contains at least one sequence that converges to zero. Let us keep in mind that some coefficients $a_{\ell,\sigma}$ in (4) may be unbounded when the parameters $(\Delta t, \Delta x)$ approach the boundary of Δ . This will restrict some of our arguments below.

When implemented on a computer, the numerical scheme (3) can of course not be used as such since it would rely on the storage of an infinite dimensional vector. We must therefore truncate the space domain, assuming for instance that the initial data are negligible outside of some domain of interest $\Omega \subset \mathbb{Z}^d$. However, as is evidenced by transport equation or by traveling waves in nonlinear equations, there is no reason why the solution to (3) should remain negligible outside of the same fixed domain $\Omega \subset \mathbb{Z}^d$ at *any* later time $n \in \mathbb{N}$. This prevents in most situations from enforcing the homogeneous Dirichlet conditions on the boundary of Ω in order to compute the restriction to the domain of interest Ω of the solution to (3). In this article, we follow a long line of research devoted to the derivation and analysis of *transparent boundary conditions*, see among many other works [AAB⁺08, AES03, BELV16, BMGN16, DZ06, EA01, ZWH08, ZE06]. We shall be concerned here with *exact* transparent boundary conditions and refer for instance to [AAB⁺08, ABS09, Ehr10, Hag99, Hal82, HY07, Sze06] for several works dedicated to the construction of approximate, more easily implementable, boundary conditions referred to as *absorbing*. In this work, we wish to understand first what are the relations satisfied by the solution to the pure Cauchy problem (3) when the data f^0, \dots, f^s vanish outside of some domain $\Omega \subset \mathbb{Z}^d$. Then among all these relations, we wish to select sufficiently many such that, when combined with the restriction of the recurrence relation (3) to Ω , we get a *stable* and hopefully *convergent* numerical initial boundary value problem. We now make our assumptions on the numerical scheme (3) precise, and then state our main results. Let us already emphasize that this article is devoted to well-posedness issues

for transparent boundary conditions only. We shall not deal here with consistency relations between the operators Q_σ in (3) and the original partial differential equation (1). Such consistency and convergence problems will be dealt with in a future work. Observe however that recurrence relations as in (3) also arise in the discretization of higher order (in time) partial differential equations such as the wave equation so our framework (3) is not restricted to the discretization of first order (in time) problems. The main technical restriction that we make here is that we focus on *scalar* problems: the coefficients $a_{\ell,\sigma}$ defining the operators Q_σ in (3) are real or complex numbers and the solutions to (3) are also real or complex valued. The extension of our work to *systems* of equations is also postponed to a future work.

1.2 Assumptions on the numerical scheme

The recurrence relation (3) is meant to be a defining equation for the sequence $u^{n+s+1} \in \ell^2(\mathbb{Z}^d)$, considering that the sequences u^n, \dots, u^{n+s} are known and belong to $\ell^2(\mathbb{Z}^d)$ (and therefore (3) is meant to define uniquely all the u^n 's, $n \geq s+1$, in terms of the initial data f^0, \dots, f^s). For future use, we not only assume that Q_{s+1} is an isomorphism on $\ell^2(\mathbb{Z}^d)$ but make the following slightly stronger assumption.

Assumption 1 (Solvability of (3)). *For all discretization parameters $(\Delta t, \Delta x) \in \Delta$, the operator Q_{s+1} is an isomorphism on $\ell^2(\mathbb{Z}^d)$, or equivalently¹:*

$$\forall \kappa \in (\mathbb{S}^1)^d, \quad \widehat{Q_{s+1}}(\kappa) := \sum_{\ell=-r}^p a_{\ell,s+1}(\Delta t, \Delta x) \kappa^\ell \neq 0, \quad \kappa^\ell := \kappa_1^{\ell_1} \cdots \kappa_d^{\ell_d}.$$

Moreover, for all $\eta = (\eta_2, \dots, \eta_d) \in \mathbb{R}^{d-1}$, Q_{s+1} satisfies the index condition:

$$\frac{1}{2i\pi} \int_{\mathbb{S}^1} \frac{\partial_{\kappa_1} \widehat{Q_{s+1}}(\kappa_1, e^{i\eta_2}, \dots, e^{i\eta_d})}{\widehat{Q_{s+1}}(\kappa_1, e^{i\eta_2}, \dots, e^{i\eta_d})} d\kappa_1 = 0. \quad (5)$$

The index condition (5) appears in many works devoted to the well-posedness of *implicit* discretizations of partial differential equations, see for instance [Str64, Osh72]. It originates in the theory of Toeplitz operators for which we refer to [GF74, Nik02] or [Lax02, Chapter 27]. Let us observe that the index condition (5) is not necessary for solving (3) on the whole space \mathbb{Z}^d . However, it will play a crucial role when we study well-posedness of the recurrence relation (3) on a ‘half-space’ $\mathbb{N} \times \mathbb{Z}^{d-1}$ or on a strip $[0; N] \times \mathbb{Z}^{d-1}$ in conjunction with the transparent boundary conditions derived below. At last, let us observe that Assumption 1 is trivially satisfied for *explicit* schemes, that is, when Q_{s+1} is the identity.

We now make a crucial assumption on the *stability* of the numerical scheme (3).

Assumption 2 (Stability of (3)). *For all $(\Delta t, \Delta x) \in \Delta$, there exists a constant $C(\Delta t, \Delta x) > 0$ such that for all initial data $f^0, \dots, f^s \in \ell^2(\mathbb{Z}^d)$, the solution to (3) satisfies the uniform bound in time*

$$\sup_{n \in \mathbb{N}} \|u^n\|_{-\infty, +\infty}^2 \leq C(\Delta t, \Delta x) \sum_{\sigma=0}^s \|f^\sigma\|_{-\infty, +\infty}^2. \quad (6)$$

¹The equivalence between the properties that Q_{s+1} is an isomorphism and that its *symbol* $\widehat{Q_{s+1}}$ does not vanish on $(\mathbb{S}^1)^d$ is based on the fact that we deal with a scalar problem. There are several other places where this restriction plays a role, though we shall not always point it out. For systems, one would need to consider the determinant of the matrix valued symbol $\widehat{Q_{s+1}}$ to characterize invertibility of Q_{s+1} .

By Fourier analysis, see [RM67, VB82, GKO95], it is a well-known fact that Assumption 2 can be completely characterized by the fulfillment of a uniform power boundedness property for the *amplification matrix* associated with (3). For future use, we therefore introduce the amplification matrix:

$$\mathcal{A}(\boldsymbol{\kappa}) := \begin{pmatrix} -\widehat{Q}_s(\boldsymbol{\kappa})/\widehat{Q}_{s+1}(\boldsymbol{\kappa}) & \dots & \dots & -\widehat{Q}_0(\boldsymbol{\kappa})/\widehat{Q}_{s+1}(\boldsymbol{\kappa}) \\ 1 & 0 & \dots & 0 \\ & \ddots & \ddots & \vdots \\ 0 & & 1 & 0 \end{pmatrix} \in \mathcal{M}_{s+1}(\mathbb{C}), \quad (7)$$

where the definition of the symbol \widehat{Q}_σ for any σ is identical to that of \widehat{Q}_{s+1} in Assumption 1, that is:

$$\forall \sigma = 0, \dots, s, \quad \widehat{Q}_\sigma(\boldsymbol{\kappa}) := \sum_{\ell=-r}^p a_{\ell,\sigma}(\Delta t, \Delta x) \boldsymbol{\kappa}^\ell.$$

In particular, we do not necessarily restrict the definition of \widehat{Q}_σ to the set $(\mathbb{S}^1)^d$. The above definition of \widehat{Q}_σ makes sense on $(\mathbb{C} \setminus \{0\})^d$. However, the amplification matrix \mathcal{A} is defined after dividing by \widehat{Q}_{s+1} , and this is possible by Assumption 1 on a neighborhood of $(\mathbb{S}^1)^d$ in \mathbb{C}^d (in what follows, we shall consider situations in which \widehat{Q}_{s+1} may vanish in $\mathbb{C} \times (\mathbb{S}^1)^{d-1}$).

Let us clarify the link between Assumption 2 and the uniform power boundedness property for \mathcal{A} since this will play a major role in the analysis below.

Lemma 1. *Let the operators Q_0, \dots, Q_{s+1} satisfy Assumptions 1 and 2. Then the amplification matrix \mathcal{A} in (7) satisfies*

$$\forall n \in \mathbb{N}, \quad \forall \boldsymbol{\kappa} \in (\mathbb{S}^1)^d, \quad \|\mathcal{A}(\boldsymbol{\kappa})^n\|^2 \leq (s+1) C(\Delta t, \Delta x), \quad (8)$$

where $C(\Delta t, \Delta x)$ is the same constant as in (6). In particular, for all $\boldsymbol{\xi} \in \mathbb{R}^d$, the dispersion relation

$$\sum_{\sigma=0}^{s+1} \widehat{Q}_\sigma(e^{i\xi_1}, \dots, e^{i\xi_d}) z^\sigma = 0, \quad (9)$$

has $s+1$ roots z_1, \dots, z_{s+1} in $\overline{\mathbb{D}}$ and those roots located on \mathbb{S}^1 are simple.

For all $(\Delta t, \Delta x) \in \boldsymbol{\Delta}$, the operators Q_σ are thus geometrically regular in the sense of [Cou09], see also [Cou13, Definition 3]. In other words the amplification matrix \mathcal{A} satisfies the uniform power boundedness (8) and if $\underline{\boldsymbol{\kappa}} \in (\mathbb{S}^1)^d$ is such that there exists $\underline{z} \in \mathbb{S}^1 \cap \text{sp}(\mathcal{A}(\underline{\boldsymbol{\kappa}}))$, then there exists a (unique) holomorphic function λ on a neighborhood \mathcal{W} of $\underline{\boldsymbol{\kappa}}$ in \mathbb{C}^d , such that

$$\begin{aligned} \lambda(\underline{\boldsymbol{\kappa}}) &= \underline{z}, \\ \det(zI + \mathcal{A}(\boldsymbol{\kappa})) &= \vartheta(\boldsymbol{\kappa}, z) (z + \lambda(\boldsymbol{\kappa})), \quad \forall (z, \boldsymbol{\kappa}) \in \mathbb{C} \times \mathcal{W}, \end{aligned}$$

with ϑ a holomorphic function of $(\boldsymbol{\kappa}, z)$ on $\mathcal{W} \times \mathbb{C}$ such that $\vartheta(\underline{\boldsymbol{\kappa}}, \underline{z}) \neq 0$, and there exists a vector valued holomorphic function $E(\boldsymbol{\kappa}) \in \mathbb{C}^{s+1}$ on \mathcal{W} that satisfies

$$\begin{aligned} E(\underline{\boldsymbol{\kappa}}) &\neq 0, \\ \forall \boldsymbol{\kappa} \in \mathcal{W}, \quad \mathcal{A}(\boldsymbol{\kappa}) E(\boldsymbol{\kappa}) &= \lambda(\boldsymbol{\kappa}) E(\boldsymbol{\kappa}). \end{aligned}$$

The interested reader will observe that the situation encountered here is a particular case of the more general geometric regularity condition of [Cou09] where several eigenvalues may ‘cross’ at $\boldsymbol{\kappa} = \underline{\boldsymbol{\kappa}}$. Here the form of the companion matrix \mathcal{A} makes such a crossing compatible with the uniform power boundedness (8) only if the crossing takes place inside the unit disk \mathbb{D} (and not on the boundary \mathbb{S}^1).

Proof of Lemma 1. The proof of the uniform bound (8) is given in [Cou13, Chapter 2] for one-dimensional explicit schemes ($d = 1$, $Q_{s+1} = I$). The extension to multidimensional implicit schemes is rather straightforward so we omit it. The dispersion relation (9) is the characteristic polynomial of \mathcal{A} , which implies that the zeroes of (9) are located in $\overline{\mathbb{D}}$. Furthermore, if one such root belongs to \mathbb{S}^1 , then the multiplicity of z as a root of (9) must equal the dimension of the eigenspace of $\mathcal{A}(\kappa)$ associated with z . Since eigenspaces of companion matrices have dimension 1, z is a simple root of (9). In particular, the corresponding eigenvalue of \mathcal{A} depends locally holomorphically on κ and one can determine an eigenvector that also depends locally holomorphically on κ . As noted in [Cou13, Lemma 7], the geometric regularity of the operators Q_σ follows automatically from the stability Assumption 2 in the scalar case. \square

The uniform bound in time we require on (3) is the (slightly relaxed) analogue of the L^2 decay property satisfied by the solutions to the partial differential equation (1). For numerical schemes with only one time level, that is when $s = 0$, then (9) is a first degree polynomial equation in z whose only root is

$$-\widehat{Q_0}(e^{i\xi_1}, \dots, e^{i\xi_d}) / \widehat{Q_1}(e^{i\xi_1}, \dots, e^{i\xi_d}),$$

and in that case, stability (in the sense of the fulfillment of (6)) is equivalent to the fact that solutions to the recurrence formula

$$Q_1 u^{n+1} + Q_0 u^n = 0,$$

satisfy the ℓ^2 decay property

$$\forall n \in \mathbb{N}, \quad \|u^{n+1}\|_{-\infty, +\infty} \leq \|u^n\|_{-\infty, +\infty}.$$

In that case, the constant $C(\Delta t, \Delta x)$ does not even depend on $(\Delta t, \Delta x)$ and can be chosen to be 1. However, when considering schemes with several time levels, there is no obvious generalization of this ℓ^2 decay since the best one can hope for is to have an energy functional that is equivalent to the norm on $\ell^2(\mathbb{Z}^d)^{s+1}$ and that is nonincreasing for solutions to (3), see [SW97]. But such an energy will strongly depend on the numerical scheme under consideration, and we therefore state Assumption 2 in this rather simplified form which encompasses a wide class of difference schemes.

What is important here is that we require the iteration (3) to satisfy a *uniform* bound in time, though we allow the constant $C(\Delta t, \Delta x)$ to depend ‘badly’ on the time and space steps. For instance we allow $C(\Delta t, \Delta x)$ to be of the form $1/\Delta t$, which would be terrible for convergence purposes, but our framework excludes numerical schemes that only satisfy a ‘Lax-Richtmyer’ bound of the form

$$\|u^{n+1}\|_{-\infty, +\infty} \leq (1 + C \Delta t) \|u^n\|_{-\infty, +\infty},$$

which in the end gives rise to an exponential bound in time for the ℓ^2 norm. Of course, the verification of Assumption 2 may restrict the set of discretization parameters Δ , and this is one reason for introducing such a set rather than considering the most general case $\Delta = (0, 1]^{d+1}$. Assumption 2 also rules out incorporating ‘lower order terms’ in the numerical scheme (3). We restrict in some sense to the principal part of (3) (just like the L^2 decay property for (1) is not invariant under lower order perturbations).

Our last main assumption appears in several anterior works on fully discrete initial boundary value problems for hyperbolic equations, see among other works [Gol77, GT81, Kre68, GKS72, Osh69]. It is also used in many *explicit* computations for deriving and analyzing discrete transparent boundary conditions for various prototype equations, see for instance [EA01, AAB⁺08, ZE06, Ehr10, DZ06, ZWH08, BELV16]. We state two possible versions of this assumption in order to highlight when the strong form of Assumption 4 is necessary and when one can only use the weak form of Assumption 3. In the defining equation (10) below, we use the decomposition $j = (j_1, j') \in \mathbb{Z} \times \mathbb{Z}^{d-1}$ for any integer $j \in \mathbb{Z}^d$. In particular r' stands for (r_2, \dots, r_d) and so on.

Assumption 3 (Noncharacteristic discrete boundary (weak form)). For $\ell_1 = -r_1, \dots, p_1$, $z \in \mathbb{C}$ and $\boldsymbol{\eta} \in \mathbb{R}^{d-1}$, let us define

$$a_{\ell_1}(z, \boldsymbol{\eta}) := \sum_{\sigma=0}^{s+1} z^\sigma \sum_{\ell'=-r'}^{p'} a_{(\ell_1, \ell'), \sigma}(\Delta t, \Delta x) e^{i \ell' \cdot \boldsymbol{\eta}}. \quad (10)$$

Then a_{-r_1} and a_{p_1} do not vanish on $\mathcal{U} \times \mathbb{R}^{d-1}$.

Assumption 4 (Noncharacteristic discrete boundary (strong form)). The functions a_{-r_1} and a_{p_1} defined in (10) do not vanish on $\overline{\mathcal{U}} \times \mathbb{R}^{d-1}$.

In what follows, we shall always consider numerical schemes that satisfy Assumption 3. This is sufficient for deriving the transparent boundary conditions for the scheme (3). However, in the stability analysis of transparent boundary conditions, it is convenient to allow $z \in \mathcal{U}$ to be arbitrarily close to the unit circle \mathbb{S}^1 . This is the reason why for showing our main stability result (Theorem 2 below), we shall make the stronger Assumption 4. The rather inconvenient feature of Assumption 4 is that it excludes numerical schemes that are first based on a semi-discretization in space of (1), giving a system of differential equations of the form

$$\frac{du_j}{dt} = \sum_{\ell=-r}^p \tilde{a}_\ell(\Delta x) u_{j+\ell},$$

and then on the application of the Crank-Nicolson rule, giving rise to the numerical scheme

$$\frac{1}{\Delta t} (u_j^{n+1} - u_j^n) = \frac{1}{2} \left(\sum_{\ell=-r}^p \tilde{a}_\ell(\Delta x) u_{j+\ell}^{n+1} + \sum_{\ell=-r}^p \tilde{a}_\ell(\Delta x) u_{j+\ell}^n \right).$$

In that case, Assumption 4 is not satisfied because either a_{-r_1} or a_{p_1} (or even both if r_1 and p_1 are nonzero) vanishes at $z = -1$. However, Assumption 4 is commonly satisfied in the discretization of hyperbolic equations by explicit methods and by some implicit schemes for parabolic equations too, see Section 5 for some examples.

We make one last, mostly technical, assumption, which restricts a little the class of numerical schemes that we consider, but that we might be able to relax later on, to the price of some future refinements in the proofs below. Some remarks on Assumption 5 are made later on in the text when we use it.

Assumption 5 (Technical restriction on the scheme). One of the following two conditions holds:

(i) Q_{s+1} is the identity operator (the scheme (3) is explicit), and for all $\boldsymbol{\eta} \in \mathbb{R}^{d-1}$ there holds

$$\begin{cases} \sum_{\ell'=-r'}^{p'} a_{(-r_1, \ell'), s}(\Delta t, \Delta x) e^{i \ell' \cdot \boldsymbol{\eta}} \neq 0, & \text{if } r_1 > 0, \\ \sum_{\ell'=-r'}^{p'} a_{(p_1, \ell'), s}(\Delta t, \Delta x) e^{i \ell' \cdot \boldsymbol{\eta}} \neq 0, & \text{if } p_1 > 0. \end{cases}$$

(ii) Q_{s+1} is not the identity (the scheme (3) is implicit), and for all $\boldsymbol{\eta} \in \mathbb{R}^{d-1}$, there holds

$$\sum_{\ell'=-r'}^{p'} a_{(-r_1, \ell'), s+1}(\Delta t, \Delta x) e^{i \ell' \cdot \boldsymbol{\eta}} \neq 0, \quad \sum_{\ell'=-r'}^{p'} a_{(p_1, \ell'), s+1}(\Delta t, \Delta x) e^{i \ell' \cdot \boldsymbol{\eta}} \neq 0.$$

At this stage, the above assumptions incorporate many standard finite difference discretizations of hyperbolic, parabolic and dispersive partial differential equations. Several examples are discussed in Section 5. Our first goal in the following paragraph is to show that the derivation of transparent boundary conditions is indeed independent of the nature of the underlying partial differential equation, but only relies on the properties encoded in Assumptions 1, 2 and 3 (we emphasize again that Assumption 5 aims mainly at simplifying one of the arguments below but should not be necessary in the most general possible framework). We shall then study the stability of the numerical scheme combining (3) with transparent boundary conditions.

1.3 Derivation of transparent numerical boundary conditions

In this paper, we focus on the derivation and analysis of transparent numerical boundary conditions when the space domain \mathbb{Z}^d is ‘truncated’ only on one side. Of course, this is still far from a realistic implementation in a computer, but our goal is to highlight, in the case of one single boundary, the main features of (3) that have an impact on the stability of the transparent numerical boundary conditions we construct below. The first step of the analysis is therefore to consider the numerical scheme (3) with initial data that vanish on a ‘half-space’, and to derive the relations satisfied by the solution to (3) in that case. This part of the analysis only uses the weak form of the assumption that the discrete boundary is noncharacteristic (Assumption 3 rather than Assumption 4). Our first main result reads as follows.

Theorem 1. *Let Assumptions 1, 2, 3 and 5 be satisfied. Then there exists a sequence $(\mathbf{\Pi}_n)_{n \in \mathbb{N}}$ of bounded operators² on $\ell^2(\mathbb{Z}^{d-1}; \mathbb{C}^{p_1+r_1})$ that satisfies*

$$\begin{aligned} \forall \delta > 0, \quad \sum_{n \in \mathbb{N}} \frac{1}{(1+\delta)^n} \|\mathbf{\Pi}_n\|_{\mathcal{B}(\ell^2(\mathbb{Z}^{d-1}))} < +\infty, \quad (\text{growth condition}), \\ \forall n \in \mathbb{N}, \quad \mathbf{\Pi}_n = \sum_{m=0}^n \mathbf{\Pi}_m \mathbf{\Pi}_{n-m}, \quad (\text{algebraic constraints}), \end{aligned}$$

and such that for all initial data $f^0, \dots, f^s \in \ell^2(\mathbb{Z}^d)$ in (3) verifying

$$\forall \sigma = 0, \dots, s, \quad \forall j_1 \leq p_1, \quad f_{(j_1, \cdot)}^\sigma = 0,$$

then the solution to (3) satisfies

$$\forall n \in \mathbb{N}, \quad \forall j_1 \leq 0, \quad \sum_{m=0}^n \mathbf{\Pi}_{n-m} \begin{pmatrix} u_{(j_1+p_1, \cdot)}^m \\ \vdots \\ u_{(j_1+1-r_1, \cdot)}^m \end{pmatrix} = 0. \quad (11)$$

When the scheme is explicit (case (i) in Assumption 5), the operator $\mathbf{\Pi}_0$ is given by

$$\mathbf{\Pi}_0 = \begin{pmatrix} 0 & 0 \\ 0 & I_{r_1} \end{pmatrix}.$$

²In the case $d = 1$, the $\mathbf{\Pi}_n$ ’s are just square matrices of size $p_1 + r_1$ since our definition of $\ell^2(\mathbb{Z}^{d-1}; \mathbb{C}^{p_1+r_1})$ then reduces to $\mathbb{C}^{p_1+r_1}$.

Let us observe that $\mathbf{\Pi}_0$ is a projector because of the algebraic constraints written for $n = 0$. When the scheme is implicit (case (ii) in Assumption 5), the operator $\mathbf{\Pi}_0$ does not have such a nice expression as in the explicit case. In particular, $\mathbf{\Pi}_0$ is most of the time a genuine *nonlocal* operator in the tangential variable $j' \in \mathbb{Z}^{d-1}$, meaning that the value of the sequence $\mathbf{\Pi}_0 v$ at an index $j' \in \mathbb{Z}^{d-1}$ does not depend on finitely many values of v close to j' but on the whole sequence v . (Of course, in the one-dimensional case $d = 1$, the $\mathbf{\Pi}_n$'s are matrices of size $p_1 + r_1$ and the discussion on non-locality becomes irrelevant.) It should be noted that in the explicit case, the $\mathbf{\Pi}_n$'s for $n \geq 1$ are also nonlocal (unless some specific - though unlikely - cancellation appears). Such nonlocality is a very common feature of transparent boundary conditions for partial differential equations in several space dimensions, see e.g. [AB01, Sze04] for the case of the Schrödinger equation.

For $n = 0, \dots, s$, the relations in (11) are trivially satisfied because the initial data f^0, \dots, f^s vanish for $j_1 \leq p_1$. Hence the relations in (11) start to be really meaningful for $n \geq s + 1$ but for reasons that will arise later on, it is useful to consider the whole set of relations (11) indexed by $n \in \mathbb{N}$ and not only by $n \geq s + 1$.

Theorem 1 provides with the set of relations (11) that are satisfied by the solution to (3) when the initial data have support in $\{j_1 > p_1\}$. In what follows, we are concerned with the numerical scheme that is obtained by combining the iteration (3) on a half-space, that is, we truncate the space domain in one spatial direction, in conjunction with discrete transparent boundary conditions obtained from (11). More precisely, we are going to consider the following numerical scheme:

$$\begin{cases} \sum_{\sigma=0}^{s+1} Q_\sigma u_j^{n+\sigma} = \Delta t F_j^{n+s+1}, & n \geq 0, \quad j_1 \geq 1, \\ \sum_{m=0}^{n+s+1} \mathbf{\Pi}_{n+s+1-m} \begin{pmatrix} u_{(p_1, \cdot)}^m \\ \vdots \\ u_{(1-r_1, \cdot)}^m \end{pmatrix} = g^{n+s+1}, & n \geq 0, \\ (u_j^0, \dots, u_j^s) = (f_j^0, \dots, f_j^s), & j_1 \geq 1 - r_1, \end{cases} \quad (12)$$

where the $\mathbf{\Pi}_n$'s are the tangential operators given by Theorem 1 and whose precise definition is given (on the Fourier side) by (33).

The discrete initial boundary value problem (12) is meant to describe, for zero interior and boundary source terms F and g , the dynamics of (3) when restricted to the half-space $\{j_1 \geq 1 - r_1\}$. There is however a little discrepancy because in Theorem 1 we have considered initial data that vanish for $j_1 \leq p_1$, while, due to the stencil of the operators Q_σ , we consider in (12) a solution (u_j^n) that is indexed by $j_1 \geq 1 - r_1$. In particular, the initial data f^0, \dots, f^s in (12) need not vanish for $j_1 \leq p_1$. We also consider the possibility of nonzero interior and boundary forcing terms F and g in view of proving stability estimates that will later be useful for convergence purposes.

The first obvious question when considering (12) is to determine whether there exists a unique solution u for given source terms F, g, f (respectively interior, boundary forcing terms, and initial data). It turns out that existence and uniqueness of a solution to (12) is *automatic* in the framework of Theorem 1. There is however a little price to pay because of the formulation of the numerical boundary conditions. Since $\mathbf{\Pi}_0$ is a projector that is not the identity (unless $p_1 = 0$), it cannot be an isomorphism on $\ell^2(\mathbb{Z}^{d-1})$ and there must necessarily be algebraic constraints on the source terms $(g^n)_{n \geq s+1}$ in (12) for proving existence of a solution to (12). These constraints are made clear in the following result³.

³Lemma 2 is purely algebraic and does not require the vector space E to be a Banach or Hilbert space, nor the linear operators to be bounded.

Lemma 2. *Let E be a vector space and let $(P_n)_{n \in \mathbb{N}}$ be a sequence of linear operators on E such that*

$$\forall n \in \mathbb{N}, \quad P_n = \sum_{m=0}^n P_m P_{n-m}.$$

Let $(y_n)_{n \in \mathbb{N}}$ be a sequence with values in E . Then there exists a sequence $(x_n)_{n \in \mathbb{N}}$ with values in E that satisfies

$$\forall n \in \mathbb{N}, \quad \sum_{m=0}^n P_{n-m} x_m = y_n, \quad (13)$$

if and only if there holds

$$\forall n \in \mathbb{N}, \quad y_n = \sum_{m=0}^n P_{n-m} y_m. \quad (14)$$

Enforcing algebraic constraints as in (14) for the source terms in (12), the unique solvability of (12) can be stated as follows.

Proposition 1. *Let Assumptions 1, 2, 3 and 5 be satisfied. Let $f^0, \dots, f^s \in \ell^2$ be the initial data for (12), and let $(g^n)_{n \geq s+1}$ be a sequence in $\ell^2(\mathbb{Z}^{d-1}; \mathbb{C}^{p_1+r_1})$ of boundary source terms for (12). With the tangential operators $(\Pi_n)_{n \geq 0}$ given in Theorem 1, let us define⁴ for $n = 0, \dots, s$, :*

$$g^n := \sum_{m=0}^n \Pi_{n-m} \begin{pmatrix} f_{(p_1, \cdot)}^m \\ \vdots \\ f_{(1-r_1, \cdot)}^m \end{pmatrix},$$

and let us further assume that the following compatibility conditions are satisfied:

$$\forall n \geq 0, \quad g^n = \sum_{m=0}^n \Pi_{n-m} g^m. \quad (15)$$

Then for all sequence of interior source terms $(F^n)_{n \geq s+1}$ with values in ℓ^2 , there exists a unique sequence $(u^n)_{n \in \mathbb{N}}$ with values in ℓ^2 solution to (12).

One drawback of Theorem 1 is that the family of operators $(\Pi_n)_{n \geq 0}$ is not uniquely defined, even when enforcing the growth condition and the algebraic constraints. There is however one choice that seems to be more natural than others, and this is the one we make in the defining equation (33). Other formulations of the transparent boundary conditions are proposed in Section 2, some of which being analogous to the one encoded in (12), and others being closer to those used in [AES03, AAB⁺08, ZE06, EA01, Ehr10, DZ06, ZWH08, BELV16] etc. We discuss how one can pass from one formulation to the other depending on the space dimension d .

Once we know that (12) has a unique solution, there remains to determine whether this solution depends continuously on the data. This is a *stability* problem, and the last requirement for Hadamard well-posedness of (12). Stability is to be understood as proving that a certain norm of the solution u to (12) can be estimated in terms of some appropriate norms of the source terms in (12). This is where the nature of the underlying partial differential equation (1) comes back into play since what we have in mind is proving an estimate for (12) that is compatible with the ‘continuous’ limit $\Delta t, \Delta x \rightarrow 0$. However, since

⁴If the initial data f^0, \dots, f^s for (12) vanish for $j_1 \leq p_1$, then g^0, \dots, g^s vanish.

the scale invariance properties of the transport, heat, Schrödinger or Airy equations are widely different, it seems hopeless at this point to encompass all possible applications within the same framework. In what follows we explore one possible notion of stability that is related to the theory of *hyperbolic* initial boundary value problems. This means that the underlying partial differential equation (1) we have in mind is a transport equation (with $\mathbf{a} \in \mathbb{R}^d$ a fixed vector):

$$\partial_t v + \sum_{j=1}^d \mathbf{a}_j \partial_j v = 0.$$

Dispersive equations such as the Schrödinger or Airy equations will be addressed in a near future with stability estimates compatible with the ones discussed in [Aud12] for the continuous problem.

1.4 Characterization of strong stability

The notion of stability that we discuss is inspired from [GKS72] and is rather restrictive in the sense that it requires controlling the trace of the solution to (12) with possibly non-homogeneous boundary forcing terms g . From now on, we consider that the ratios $\Delta t / \Delta x_i$, $i = 1, \dots, d$, are constant, which means that the set Δ of discretization parameters is a semi-open interval

$$\Delta = \{(\Delta t, \Delta t / \lambda_1, \dots, \Delta t / \lambda_d), \quad \Delta t \in (0, 1]\},$$

where $\lambda_1, \dots, \lambda_d$ are fixed positive numbers. We also assume that the coefficients $a_{\ell, \sigma}$ in (4) only depend on the time and space steps through the ratios $\lambda_1, \dots, \lambda_d$. In other words, keeping Δt as the only free small parameter, we assume that the operators Q_σ are *independent* of Δt , which means that the scheme (12) is also independent of Δt (since one can rename the interior source term $\Delta t F_j^{n+s+1}$ fictitiously as \tilde{F}_j^{n+s+1}). The scaling invariance $\Delta t / \Delta x_i = \text{cst}$ is reminiscent of the scaling invariance $(t, x) \rightarrow (\alpha t, \alpha x)$ of the underlying hyperbolic equation we have in mind. The estimate (16) below is also the direct analogue of the weighted in time estimates discussed in [BGS07, Chapter 4] and that are also invariant by the scaling $(t, x) \rightarrow (\alpha t, \alpha x)$ (the Laplace parameter γ below being then rescaled as $\gamma \rightarrow \gamma / \alpha$). As already mentioned, we plan to adapt the continuous *dispersive* estimates of [Aud12] to the framework of finite difference schemes and transparent boundary conditions in a near future. We now introduce the following terminology.

Definition 1 (Strong stability [GKS72]). *The finite difference approximation (12) is said to be strongly stable if there exists a constant C such that for all $\gamma > 0$ and all $\Delta t \in (0, 1]$, the solution⁵ (u_j^n) to (12) with $(f_j^0) = \dots = (f_j^s) = 0$ satisfies the estimate:*

$$\begin{aligned} & \frac{\gamma}{\gamma \Delta t + 1} \sum_{n \geq s+1} \Delta t e^{-2\gamma n \Delta t} \|u^n\|_{1-r_1, +\infty}^2 + \sum_{n \geq s+1} \sum_{j_1=1-r_1}^{p_1} \Delta t e^{-2\gamma n \Delta t} \|u_{(j_1, \cdot)}^n\|_{\ell^2(\mathbb{Z}^{d-1})}^2 \\ & \leq C_1 \left\{ \frac{\gamma \Delta t + 1}{\gamma} \sum_{n \geq s+1} \Delta t e^{-2\gamma n \Delta t} \|F^n\|_{1, +\infty}^2 + \sum_{n \geq s+1} \Delta t e^{-2\gamma n \Delta t} \|g^n\|_{\ell^2(\mathbb{Z}^{d-1})}^2 \right\}. \quad (16) \end{aligned}$$

⁵Here we tacitly assume that the boundary forcing terms and vanishing initial data satisfy the compatibility conditions (15) so that a solution to (12) does exist.

The main point in the estimate (16) is that the $\ell_{n,j'}^2$ norm of the trace of the solution is estimated with the same ‘weight’ as γ times the $\ell_{n,j}^2$ norm. Moreover, there is no loss of ‘derivative’ from the source terms in (12) to the solution in the estimate (16). We exhibit below a necessary and sufficient condition for (12) to be strongly stable. The analysis is inspired from and bears quite some resemblance with [Cou15a], see also [Tre84]. Before stating our main result, we introduce some more terminology.

Definition 2 (Non-glancing scheme). *Let the operators Q_0, \dots, Q_{s+1} satisfy Assumptions 1 and 2. The scheme (3) is said to be non-glancing if for all $\underline{\kappa} \in (\mathbb{S}^1)^d$ such that there exists $\underline{z} \in \mathbb{S}^1 \cap \text{sp}(\mathcal{A}(\underline{\kappa}))$, then the holomorphic eigenvalue λ of \mathcal{A} given in Lemma 2 satisfies*

$$\frac{\partial \lambda}{\partial \kappa_1}(\underline{\kappa}) \neq 0.$$

As discussed in [Cou15a], several standard discretizations of transport equations such as the upwind, Lax-Friedrichs and Lax-Wendroff schemes are non-glancing. At the opposite, and as already noticed in [Tre84], the leap-frog and Crank-Nicolson approximations of the transport equation admit glancing wave packets (while the underlying partial differential equation does not !). Such examples are discussed in Section 5. Our main result reads as follows. We emphasize that from now on we enforce the stronger Assumption 4 rather than its weak version (Assumption 3).

Theorem 2. *Let Assumptions 1, 2, 4 and 5 be satisfied. Then the scheme (12) is strongly stable in the sense of Definition 1 if and only if the scheme (3) for the pure Cauchy problem is non-glancing.*

Assume furthermore that for all $\xi \in \mathbb{R}^d$, the roots to the dispersion relation (9) are simple. Then if the scheme (3) is non-glancing, there exists a constant $C > 0$ such that for all $\gamma > 0$ and all $\Delta t \in (0, 1]$, the solution⁶ (u_j^n) to (12) with $f^0, \dots, f^s \in \ell^2$ satisfies the estimate:

$$\begin{aligned} & \sup_{n \in \mathbb{N}} e^{-2\gamma n \Delta t} \|u^n\|_{1-r_1, +\infty}^2 + \sum_{n \geq s+1} \sum_{j_1=1-r_1}^{p_1} \Delta t e^{-2\gamma n \Delta t} \|u_{(j_1, \cdot)}^n\|_{\ell^2(\mathbb{Z}^{d-1})}^2 \\ & \leq C \left\{ \sum_{\sigma=0}^s \|f^\sigma\|_{1-r_1, +\infty}^2 + \frac{\gamma \Delta t + 1}{\gamma} \sum_{n \geq s+1} \Delta t e^{-2\gamma n \Delta t} \|F^n\|_{1, +\infty}^2 + \sum_{n \geq s+1} \Delta t e^{-2\gamma n \Delta t} \|g^n\|_{\ell^2(\mathbb{Z}^{d-1})}^2 \right\}. \end{aligned} \quad (17)$$

Let us observe that in (17), we could have added for free the quantity

$$\frac{\gamma}{\gamma \Delta t + 1} \sum_{n \geq 0} \Delta t e^{-2\gamma n \Delta t} \|u^n\|_{1-r_1, +\infty}^2,$$

on the left hand side, since it is controlled by the stronger ‘semigroup’ norm

$$\sup_{n \in \mathbb{N}} e^{-2\gamma n \Delta t} \|u^n\|_{1-r_1, +\infty}^2,$$

uniformly in γ and Δt . The main point in Theorem 2 is that for non-glancing schemes, one can control both semigroup and trace norms of the solution to (12), including for the case of non-homogeneous boundary conditions. Such strong stability estimates are relevant for two purposes: first they allow for a

⁶Here again we tacitly assume that the boundary forcing terms and the (nonzero) initial data satisfy the compatibility conditions (15) so that a solution to (12) does exist.

convergence analysis relying on suitable consistency estimates, and second such strong stability estimates persist under any small modification of the numerical boundary conditions. This means that one may be able to prove first stability and then convergence for a class of *absorbing* boundary conditions that is based on a sufficiently good approximation of the tangential operators $(\mathbf{\Pi}_n)_{n \geq 0}$. We plan to investigate this issue in a future work, for instance when the operators $\mathbf{\Pi}_n$ are approximated by the so-called sum of exponentials, see e.g. [AES03].

The plan of the paper is as follows. In Section 2, we derive the transparent numerical boundary conditions and prove Theorem 1. We also discuss alternative formulations of the transparent numerical boundary conditions and make some constructions more explicit in the special case $p_1 = r_1 = 1$ which occurs in many practical examples. In Section 3, we show solvability for the numerical scheme (12) and prove Lemma 2 and Proposition 1. We also briefly describe, for $d = 1$, the case where the space domain \mathbb{Z} is truncated on either side, leading to a *finite dimensional* problem on an interval $[0, N]$. Section 4 is devoted to the proof of Theorem 2 which characterizes strong stability in terms of non-existence of glancing wave packets. Eventually we discuss in Section 5 several examples related to transport, diffusion or dispersive equations.

2 Derivation of transparent numerical boundary conditions

In this Section, we prove Theorem 1 and construct transparent boundary conditions in the rather wide framework covered by Assumptions 1, 2, 3 and 5.

2.1 Proof of Theorem 1

The proof splits in several steps. In the first step, we compute the transparent conditions on the ‘Laplace-Fourier side’. In several references, this calculation is performed by using the so-called \mathcal{Z} -transform, which is the discrete analogue of the Laplace transform. The second step consists in going back to the original time and space variables by using the inverse Laplace-Fourier transform. This requires specific attention in order to show that the ‘causality’ principle is preserved, meaning that in (11), the relation that should eventually contribute to determining u^n does not involve some $u^{n'}$ with $n' > n$. Writing the relations (11) in a causal way in which ‘future does not affect the past’ amounts to proving that the Laurent expansions of several objects that are holomorphic on \mathcal{U} do have a limit at infinity (and therefore the Laurent expansions only involve nonpositive powers of $z \in \mathcal{U}$). Justifying that such limits at infinity exist is the purpose of the second step of the proof and this is the reason why we make the technical Assumption 5. Once this is achieved, the conclusions of Theorem 1 follow rather easily.

- Step 1. The transparent conditions on the Laplace-Fourier side.

We consider the iteration (3) with initial data in ℓ^2 vanishing for $j_1 \leq p_1$. Thanks to Assumption 2, we know that the solution to (3) satisfies

$$\forall \gamma > 0, \quad \sum_{n \geq 0} \Delta t e^{-2\gamma n \Delta t} \|u^n\|_{-\infty, +\infty}^2 < +\infty. \quad (18)$$

In particular, for all fixed $x_1 \in \mathbb{R}$, the piecewise constant function $u(\cdot, x_1, \cdot)$ in (2) has a well-defined Laplace-Fourier transform on $\mathbb{C}^+ \times \mathbb{R}^{d-1}$, where we let \mathbb{C}^+ denote the set of complex numbers with positive real part. Here the Laplace transform refers to the time variable t , and the (partial) Fourier transform refers to the tangential space variables $x' := (x_2, \dots, x_d)$. Even though the x_1 variable lies

in \mathbb{R} , we do not perform Fourier transform with respect to x_1 . We rather stick to the original discrete variable $j_1 \in \mathbb{Z}$ and feel free to use the notation $u_{j_1}(t, x')$ which is entirely analogous to the notation introduced in (2) for the step function u . The dual variables to (t, x') are denoted (τ, ξ') below, with $\tau = \gamma + i\theta$ and $\xi' = (\xi_2, \dots, \xi_d)$, and the Laplace-Fourier transform of u_{j_1} is denoted \widehat{u}_{j_1} . Moreover, given $(\tau, \xi') \in \mathbb{C}^+ \times \mathbb{R}^{d-1}$, we always use the notation

$$z := e^{\tau \Delta t} \in \mathcal{U}, \quad \boldsymbol{\eta} := (\xi_2 \Delta x_2, \dots, \xi_d \Delta x_d) \in \mathbb{R}^{d-1}.$$

The difference between the Laplace-Fourier transform used here and the \mathcal{Z} -transform used in [AES03, AAB⁺08, ZE06, Ehr10, BELV16] is only of a multiplicative - though *crucial* - factor. Indeed, we compute

$$\widehat{u}_{j_1}(\tau, \xi') = \frac{1 - z^{-1}}{\tau} \sum_{n \geq 0} z^{-n} \int_{\mathbb{R}^{d-1}} e^{-i x' \cdot \xi'} u_{j_1}^n(x') dx', \quad (19)$$

at least if, for all $n \in \mathbb{N}$, the sequence $(u_{j_1, \cdot}^n)$ belongs to $\ell^1(\mathbb{Z}^{d-1})$. Otherwise we use the continuous extension of the Fourier transform to $L^2(\mathbb{R}^{d-1})$ in case $(u_{j_1, \cdot}^n)$ belongs to $\ell^2(\mathbb{Z}^{d-1})$ without belonging to $\ell^1(\mathbb{Z}^{d-1})$.

Applying Plancherel's Theorem, we obtain from the bound (18) the property:

$$\forall \gamma > 0, \quad \sum_{j_1 \in \mathbb{Z}} \int_{\mathbb{R} \times \mathbb{R}^{d-1}} |\widehat{u}_{j_1}(\gamma + i\theta, \xi')|^2 d\theta d\xi' < +\infty.$$

In particular, for any $\gamma > 0$, the sequence $(\widehat{u}_{j_1}(\gamma + i\theta, \xi'))_{j_1 \in \mathbb{Z}}$ belongs to $\ell^2(\mathbb{Z})$ for almost every $(\theta, \xi') \in \mathbb{R} \times \mathbb{R}^{d-1}$. We omit below the possible negligible set of those (θ, ξ') for which the sequence is not in ℓ^2 , and make as if this negligible set is always empty. This has no consequence of course in the proof.

With the above notation, it is rather straightforward to derive the recurrence relation satisfied by the sequence $(\widehat{u}_{j_1}(\gamma + i\theta, \xi'))_{j_1 \in \mathbb{Z}}$. Namely, we first apply the partial Fourier transform with respect to x' and get⁷

$$\begin{cases} \sum_{\sigma=0}^{s+1} Q_{\sigma}^{\#}(\boldsymbol{\eta}) \widehat{u_{j_1}^{n+\sigma}}(\xi') = 0, & n \geq 0, \quad j_1 \in \mathbb{Z}, \\ \left(\widehat{u_{j_1}^0}, \dots, \widehat{u_{j_1}^s} \right)(\xi') = \left(\widehat{f_{j_1}^0}, \dots, \widehat{f_{j_1}^s} \right)(\xi'), & j_1 \in \mathbb{Z}, \end{cases}$$

with

$$\forall \sigma = 0, \dots, s, \quad Q_{\sigma}^{\#}(\boldsymbol{\eta}) := \sum_{\ell_1=-r_1}^{p_1} \left(\sum_{\ell'=-r'}^{p'} a_{(\ell_1, \ell'), \sigma}(\Delta t, \Delta x) e^{i \ell' \cdot \boldsymbol{\eta}} \right) \mathbf{S}_1^{\ell_1}, \quad (\mathbf{S}_1^{\ell_1} w)_{j_1} := w_{j_1 + \ell_1}. \quad (20)$$

Let us recall that $\boldsymbol{\eta}$ is a placeholder for $(\xi_2 \Delta x_2, \dots, \xi_d \Delta x_d)$. We then use the expression (19) of the Laplace-Fourier transform \widehat{u}_{j_1} to compute the relation

$$\forall j_1 \in \mathbb{Z}, \quad \sum_{\ell_1=-r_1}^{p_1} a_{\ell_1}(z, \boldsymbol{\eta}) \widehat{u_{j_1+\ell_1}}(\tau, \xi') = F_{j_1}(\tau, \xi'), \quad (21)$$

⁷It is convenient here to use also the hat notation for denoting the partial Fourier transform with respect to the tangential spatial variables x' .

where the functions a_{-r_1}, \dots, a_{p_1} are defined in (10) and the source term F_{j_1} in (21) is given by the relation

$$F_{j_1}(\tau, \xi') := \frac{1 - z^{-1}}{\tau} \sum_{m=0}^s \sum_{\sigma=m}^s z^{1+\sigma-m} Q_{1+\sigma}^\#(\eta) \widehat{f_{j_1}^m}(\xi').$$

The precise expression of F_{j_1} is not very useful. What is important is that, because of the support assumption on the initial data in Theorem 1, F_{j_1} is zero for any index $j_1 \leq 0$. Thanks to Assumption 3, we also know that both $a_{-r_1}(z, \eta)$ and $a_{p_1}(z, \eta)$ are nonzero because $\tau \in \mathbb{C}^+$ and therefore $z = e^{\tau \Delta t} \in \mathcal{U}$. Hence (21) is a recurrence relation of order (exactly equal to) $p_1 + r_1$ for the sequence $(\widehat{u_{j_1}}(\tau, \xi'))_{j_1 \in \mathbb{Z}}$. Introducing the vector

$$U_{j_1}(\tau, \xi') := \begin{pmatrix} \widehat{u_{j_1+p_1-1}}(\tau, \xi') \\ \vdots \\ \widehat{u_{j_1-r_1}}(\tau, \xi') \end{pmatrix} \in \mathbb{C}^{p_1+r_1},$$

as well as the companion matrix

$$\forall (z, \eta) \in \mathcal{U} \times \mathbb{R}^{d-1}, \quad \mathbb{M}(z, \eta) := \begin{pmatrix} -\frac{a_{p_1-1}(z, \eta)}{a_{p_1}(z, \eta)} & \dots & \dots & -\frac{a_{-r_1}(z, \eta)}{a_{p_1}(z, \eta)} \\ 1 & 0 & \dots & 0 \\ & \ddots & \ddots & \vdots \\ 0 & & 1 & 0 \end{pmatrix} \in \mathcal{M}_{p_1+r_1}(\mathbb{C}), \quad (22)$$

we can equivalently rewrite (21) as

$$U_{j_1+1}(\tau, \xi') - \mathbb{M}(z, \eta) U_{j_1}(\tau, \xi') = \frac{1}{a_{p_1}(z, \eta)} \begin{pmatrix} F_{j_1}(\tau, \xi') \\ 0 \\ \vdots \\ 0 \end{pmatrix}.$$

In particular, there holds

$$\forall j_1 \leq 0, \quad U_{j_1+1}(\tau, \xi') = \mathbb{M}(z, \eta) U_{j_1}(\tau, \xi'). \quad (23)$$

The relations (11) of Theorem 1 are direct consequences of the recurrence relation (23). Some fundamental properties of the matrix $\mathbb{M}(z, \eta)$ are stated in the following Lemma.

Lemma 3. *Let Assumptions 1, 2, 3 and 5 be satisfied. Then for all $(z, \eta) \in \mathcal{U} \times \mathbb{R}^{d-1}$, the matrix $\mathbb{M}(z, \eta)$ in (22) is invertible and has no eigenvalue on \mathbb{S}^1 . If moreover Assumption 5 is satisfied, then $\mathbb{M}(z, \eta)$ has r_1 eigenvalues (counted with multiplicity) in \mathbb{D} and the remaining p_1 eigenvalues in \mathcal{U} .*

Let us take the result of Lemma 3 for granted for a while. We can therefore introduce the eigenprojectors $\Pi^{s,u}(z, \eta)$ associated with the decomposition of $\mathbb{C}^{p_1+r_1}$ into

$$\mathbb{C}^{p_1+r_1} = \mathbb{E}^s(z, \eta) \oplus \mathbb{E}^u(z, \eta). \quad (24)$$

Here \mathbb{E}^s refers to the *stable* subspace, that is the generalized eigenspace of $\mathbb{M}(z, \eta)$ associated with the eigenvalues in \mathbb{D} and \mathbb{E}^u refers to the *unstable* subspace associated with the eigenvalues of \mathbb{M} in \mathcal{U} . By Lemma 3, the rank of $\Pi^s(z, \eta)$ is r_1 and the rank of $\Pi^u(z, \eta)$ is p_1 .

Using the recurrence relation (23) and the fact that $(U_{j_1}(\tau, \xi'))_{j_1 \in \mathbb{Z}}$ belongs to $\ell^2(\mathbb{Z})$, the vector $U_1(\tau, \xi')$ must necessarily belong to the unstable subspace $\mathbb{E}^u(z, \eta)$ and this property is propagated to

any $j_1 \leq 0$ by the recurrence (23) since $\mathbb{E}^u(z, \boldsymbol{\eta})$ is invariant by \mathbb{M} and \mathbb{M}^{-1} . Using the projectors $\Pi^{s,u}(z, \boldsymbol{\eta})$ associated with (24), we have obtained:

$$\forall j_1 \leq 1, \quad \Pi^s(z, \boldsymbol{\eta}) U_{j_1}(\tau, \boldsymbol{\xi}') = 0. \quad (25)$$

The relations (25) are the Laplace-Fourier counterparts of the relations (11) stated in Theorem 1. In what follows, we are going to prove Lemma 3. Then we shall explain how one can rewrite (25) in the original physical space, which amounts to determining the inverse Laplace-Fourier transform of the left hand side in (25).

Proof of Lemma 3. The proof is mostly the same as in [Kre68], but since we intend to show that the argument is not restricted to discretized *hyperbolic* problems, we reproduce it here for the sake of completeness. Let $(z, \boldsymbol{\eta}) \in \mathcal{U} \times \mathbb{R}^{d-1}$. The determinant of the matrix $\mathbb{M}(z, \boldsymbol{\eta})$ in (22) equals

$$(-1)^{p_1+r_1} a_{-r_1}(z, \boldsymbol{\eta}) / a_{p_1}(z, \boldsymbol{\eta}),$$

this quantity being nonzero due to Assumption 3. More generally, the characteristic polynomial of $\mathbb{M}(z, \boldsymbol{\eta})$ is

$$\frac{1}{a_{p_1}(z, \boldsymbol{\eta})} \sum_{\ell_1=-r_1}^{p_1} a_{\ell_1}(z, \boldsymbol{\eta}) X^{\ell_1+r_1},$$

and the eigenvalues of $\mathbb{M}(z, \boldsymbol{\eta})$ are the roots $\kappa_1 \in \mathbb{C} \setminus \{0\}$ to the dispersion relation

$$\sum_{\sigma=0}^{s+1} \widehat{Q}_\sigma(\kappa_1, e^{i\eta_2}, \dots, e^{i\eta_d}) z^\sigma = 0.$$

In particular, κ_1 does not belong to the unit circle \mathbb{S}^1 for otherwise (9) would have a root z in \mathcal{U} for some $\boldsymbol{\xi} \in \mathbb{R}^d$ (and this would contradict Lemma 1). This means that the eigenvalues of $\mathbb{M}(z, \boldsymbol{\eta})$ split into two groups: the *stable* ones in \mathbb{D} and the *unstable* ones in \mathcal{U} , which gives rise to the decomposition (24). It remains to determine the dimension of each of the vector spaces in (24).

We start with the case of explicit schemes, that is, case (i) in Assumption 5. Following [Kre68, Lemma 2] (see also [Cou13, Lemma 15] for a more detailed exposition), the number of stable eigenvalues is computed by analyzing their behavior when z tends to infinity. Since $\mathcal{U} \times \mathbb{R}^{d-1}$ is connected, the number of eigenvalues of $\mathbb{M}(z, \boldsymbol{\eta})$ in \mathbb{D} does not depend on $(z, \boldsymbol{\eta})$. Let now $\boldsymbol{\eta}$ be fixed. As z tends to infinity, all stable eigenvalues of $\mathbb{M}(z, \boldsymbol{\eta})$ converge to zero. This property holds for otherwise, there would exist a sequence (z_n) with $|z_n| > n$ and there would exist a sequence (κ_n) such that $\kappa_n \in \mathbb{D} \cap \text{sp}(\mathbb{M}(z_n, \boldsymbol{\eta}))$, and $\inf_n |\kappa_n| > 0$. Up to extracting and relabeling, we can assume that the sequence (κ_n) converges towards some nonzero complex number $\underline{\kappa}$. We now pass to the limit in the expression

$$\frac{1}{z_n^{s+1}} \sum_{\ell_1=-r_1}^{p_1} a_{\ell_1}(z_n, \boldsymbol{\eta}) \kappa_n^{\ell_1+r_1} = 0,$$

where the rescaling by z_n^{-s-1} has been made in order to have (recall here that we consider explicit schemes):

$$\frac{1}{z_n^{s+1}} a_{\ell_1}(z_n, \boldsymbol{\eta}) \longrightarrow \delta_{\ell_1 0},$$

see (10). Hence we get $\underline{\kappa}^{r_1} = 0$, which is a contradiction. All the stable eigenvalues of $\mathbb{M}(z, \boldsymbol{\eta})$ tend to zero as z tends to infinity (and similarly, all the unstable eigenvalues tend to infinity as z tends to infinity). Setting $z = 1/Z$, the number of stable eigenvalues of $\mathbb{M}(z, \boldsymbol{\eta})$ is computed for Z small by counting the number of roots close to zero to the equation

$$D(\kappa, Z) := Z^{s+1} \sum_{\ell_1=-r_1}^{p_1} a_{\ell_1}(1/Z, \boldsymbol{\eta}) \kappa^{\ell_1+r_1} = 0.$$

Since D is a polynomial in (κ, Z) with $D(\kappa, 0) = \kappa^{r_1}$, there are r_1 stable eigenvalues of $\mathbb{M}(z, \boldsymbol{\eta})$. The p_1 other eigenvalues of $\mathbb{M}(z, \boldsymbol{\eta})$ are the unstable ones.

We now deal with the case of implicit schemes, that is, case (ii) in Assumption 5. Once again, we analyze the behavior of the eigenvalues of $\mathbb{M}(z, \boldsymbol{\eta})$ as z tends to infinity. For fixed $\boldsymbol{\eta}$ and κ_1 , we compute

$$\lim_{z \rightarrow \infty} \frac{1}{z^{s+1}} \sum_{\ell_1=-r_1}^{p_1} a_{\ell_1}(z, \boldsymbol{\eta}) \kappa_1^{\ell_1+r_1} = \sum_{\ell_1=-r_1}^{p_1} \left(\sum_{\ell'=-r'}^{p'} a_{(\ell_1, \ell'), s+1}(\Delta t, \Delta x) e^{i \ell' \cdot \boldsymbol{\eta}} \right) \kappa_1^{\ell_1+r_1}.$$

Using Assumptions 1 and 5, we know that the equation in κ_1 :

$$\kappa_1^{r_1} \widehat{Q_{s+1}}(\kappa_1, e^{i \eta_2}, \dots, e^{i \eta_d}) = \sum_{\ell_1=-r_1}^{p_1} \left(\sum_{\ell'=-r'}^{p'} a_{(\ell_1, \ell'), s+1}(\Delta t, \Delta x) e^{i \ell' \cdot \boldsymbol{\eta}} \right) \kappa_1^{\ell_1+r_1} = 0,$$

is a polynomial equation of degree $p_1 + r_1$, that its roots are nonzero and that it has no root on \mathbb{S}^1 . Furthermore, the residue Theorem [Rud87] shows that the integral on the left hand side of (5) equals the number of zeroes in \mathbb{D} of the holomorphic function $\widehat{Q_{s+1}}(\cdot, e^{i \eta_2}, \dots, e^{i \eta_d})$ minus the number of its poles in \mathbb{D} (zeroes and poles being counted with multiplicity). By Assumption 5, we know that there is only one pole at the origin with order r_1 , so the number of zeroes in \mathbb{D} is r_1 . Summarizing, we have shown that the equation

$$\sum_{\ell_1=-r_1}^{p_1} \left(\sum_{\ell'=-r'}^{p'} a_{(\ell_1, \ell'), s+1}(\Delta t, \Delta x) e^{i \ell' \cdot \boldsymbol{\eta}} \right) \kappa_1^{\ell_1+r_1} = 0,$$

has r_1 roots in \mathbb{D} and it must therefore also have p_1 roots in \mathcal{U} . By the Rouché Theorem [Rud87], this implies that for any sufficiently large z , the equation

$$\sum_{\ell_1=-r_1}^{p_1} a_{\ell_1}(z, \boldsymbol{\eta}) \kappa_1^{\ell_1+r_1} = 0,$$

also has r_1 roots in \mathbb{D} and p_1 roots in \mathcal{U} . □

- Step 2. The limit at infinity of the stable and unstable subspaces.

Let us first observe that because of Lemma 3, the stable and unstable spaces $E^{s,u}$ in (24) depend holomorphically on $z \in \mathcal{U}$ and analytically on $\boldsymbol{\eta} \in \mathbb{R}^{d-1}$. Moreover, they are 2π -periodic with respect to each coordinate of $\boldsymbol{\eta}$ because the matrix \mathbb{M} itself is 2π -periodic with respect to each coordinate of $\boldsymbol{\eta}$, see (10). Hence the projectors $\Pi^{s,u}$, which can be defined by integrals of $(wI - \mathbb{M})^{-1}$ on suitable contours [Bau85], depend holomorphically on $z \in \mathcal{U}$ and analytically on $\boldsymbol{\eta} \in \mathbb{R}^{d-1}$. Assumption 5 comes back into

play for studying the limit of the projectors $\Pi^{s,u}$ as z tends to infinity, which amounts to determining the limits of $E^{s,u}$ as z tends to infinity.

Let us first consider case (ii) in Assumption 5, that is, the case of implicit schemes. We have already seen in the proof of Lemma 3 that for all $\boldsymbol{\eta} \in \mathbb{R}^{d-1}$, the function

$$\kappa_1 \in \mathbb{C} \setminus \{0\} \mapsto \widehat{Q_{s+1}}(\kappa_1, e^{i\eta_2}, \dots, e^{i\eta_d}),$$

has r_1 zeroes in $\mathbb{D} \setminus \{0\}$ and p_1 zeroes in \mathcal{U} (as always, zeroes are counted with multiplicity). We can then easily determine the behavior of the spectral projectors $\Pi^{s,u}(z, \boldsymbol{\eta})$ as z tends to infinity. Indeed, by the definition (22) of $\mathbb{M}(z, \boldsymbol{\eta})$, and the fact that both a_{-r_1} and a_{p_1} are polynomials of degree $s+1$ in z (Assumption 5), we find that $\mathbb{M}(z, \boldsymbol{\eta})$ has a limit as z tends to infinity. This limit is given by

$$\begin{pmatrix} -\frac{a_{\infty, p_1-1}(\boldsymbol{\eta})}{a_{\infty, p_1}(\boldsymbol{\eta})} & \dots & \dots & -\frac{a_{\infty, -r_1}(\boldsymbol{\eta})}{a_{\infty, p_1}(\boldsymbol{\eta})} \\ 1 & 0 & \dots & 0 \\ & \ddots & \ddots & \vdots \\ 0 & & 1 & 0 \end{pmatrix}, \quad a_{\infty, \ell_1}(\boldsymbol{\eta}) := \sum_{\ell'=-r'}^{p'} a_{(\ell_1, \ell'), s+1}(\Delta t, \Delta x) e^{i\ell' \cdot \boldsymbol{\eta}}, \quad (26)$$

and we know from the previous arguments that this matrix has r_1 eigenvalues in $\mathbb{D} \setminus \{0\}$ and p_1 eigenvalues in \mathcal{U} . In other words, the singularity at $Z = 0$ of the projectors $\Pi^{s,u}(1/Z, \boldsymbol{\eta})$ is removable since the splitting between stable and unstable eigenvalues persists up to $Z = 0$. In what follows, we let $\mathbb{E}^s(\infty, \boldsymbol{\eta})$, resp. $\mathbb{E}^u(\infty, \boldsymbol{\eta})$, denote the stable, resp. unstable, subspace of the matrix $\mathbb{M}(\infty, \boldsymbol{\eta})$ in (26). We also let $\Pi^{s,u}(\infty, \boldsymbol{\eta})$ denote the projectors associated with the decomposition

$$\mathbb{C}^{p_1+r_1} = \mathbb{E}^s(\infty, \boldsymbol{\eta}) \oplus \mathbb{E}^u(\infty, \boldsymbol{\eta}),$$

which holds for any $\boldsymbol{\eta} \in \mathbb{R}^{d-1}$.

Let us now turn to the case of explicit schemes, case (i) in Assumption 5. This is the case treated in [Kre68, Cou09]. From the proof of Lemma 3, we already know that the eigenvalues of $\mathbb{M}(z, \boldsymbol{\eta})$ are the roots κ_1 to the equation

$$\sum_{\ell_1=-r_1}^{p_1} \kappa_1^{\ell_1+r_1} a_{\ell_1}(z, \boldsymbol{\eta}) = 0.$$

There are r_1 stable roots, all of them converging to zero as z tends to infinity, and there are p_1 unstable roots, all of them tending to infinity as z tends to infinity. We are going to determine an asymptotic expansion of the eigenvalues as z tends to infinity. (Of course determining this behavior makes sense only if r_1 and/or p_1 is nonzero, and this is the reason why in Assumption 5 we have stated conditions only when these integers are nonzero, which we assume from now on.) We introduce the function

$$D : (Z, \kappa_1) \mapsto Z^{s+1} \sum_{\ell_1=-r_1}^{p_1} \kappa_1^{\ell_1+r_1} a_{\ell_1}(1/Z, \boldsymbol{\eta}),$$

that is a polynomial function of (Z, κ_1) , and that satisfies

$$D(0, \kappa_1) = \kappa_1^{r_1}, \quad \frac{\partial D}{\partial Z}(0, 0) = \lim_{Z \rightarrow 0} Z^s a_{-r_1}(1/Z, \boldsymbol{\eta}) = \sum_{\ell'=-r'}^{p'} a_{(-r_1, \ell'), s}(\Delta t, \Delta x) e^{i\ell' \cdot \boldsymbol{\eta}} \neq 0.$$

Applying the Puiseux expansions theory, for which we refer to [Bau85], the r_1 eigenvalues of $\mathbb{M}(1/Z, \boldsymbol{\eta})$ close to zero thus have an asymptotic expansion of the form

$$\kappa_1 \sim c Z^{1/r_1} + O(Z^{2/r_1}), \quad c \neq 0.$$

Recall that the frequency $\boldsymbol{\eta}$ is fixed here. In particular, the r_1 stable eigenvalues of $\mathbb{M}(z, \boldsymbol{\eta})$ are simple for all sufficiently large z , $\boldsymbol{\eta}$ being fixed (the splitting comes from the r_1 possible branches for the r_1 -th root of $1/z$). In a similar way, we can prove that the p_1 unstable eigenvalues of $\mathbb{M}(z, \boldsymbol{\eta})$ have the asymptotic expansion

$$\kappa_1 \sim d z^{1/p_1} + O(1), \quad d \neq 0.$$

At this stage, we know that for all sufficiently large z , the eigenvalues of $\mathbb{M}(z, \boldsymbol{\eta})$ are simple; the stable ones behave like z^{-1/r_1} and the unstable ones behave like z^{1/p_1} as z tends to infinity. It remains to compute the limit of $\Pi^s(z, \boldsymbol{\eta})$ (the limit of Π^u is directly obtained by using $\Pi^u = I - \Pi^s$), which will also give the limit of $\mathbb{E}^{s,u}(z, \boldsymbol{\eta})$ as z tends to infinity.

Using the expression of the eigenvectors of a companion matrix, we know that the Vandermonde matrix

$$\begin{pmatrix} \kappa_{1,1}^{p_1+r_1-1} & \cdots & \kappa_{1,p_1+r_1}^{p_1+r_1-1} \\ \vdots & & \vdots \\ 1 & \cdots & 1 \end{pmatrix},$$

diagonalizes \mathbb{M} for z large enough (so that all eigenvalues are simple). We label the eigenvalues in such a way that $\kappa_{1,1}, \dots, \kappa_{1,r_1}$ are the stable ones, and $\kappa_{1,r_1+1}, \dots, \kappa_{1,r_1+p_1}$ are the unstable ones. With this convention, there holds

$$\Pi^s(z, \boldsymbol{\eta}) = \begin{pmatrix} \kappa_{1,1}^{p_1+r_1-1} & \cdots & \kappa_{1,p_1+r_1}^{p_1+r_1-1} \\ \vdots & & \vdots \\ 1 & \cdots & 1 \end{pmatrix} \begin{pmatrix} I_{r_1} & 0 \\ 0 & 0 \end{pmatrix} \begin{pmatrix} \kappa_{1,1}^{p_1+r_1-1} & \cdots & \kappa_{1,p_1+r_1}^{p_1+r_1-1} \\ \vdots & & \vdots \\ 1 & \cdots & 1 \end{pmatrix}^{-1}.$$

Our aim is to show the property

$$\lim_{z \rightarrow \infty} \Pi^s(z, \boldsymbol{\eta}) = \begin{pmatrix} 0 & 0 \\ 0 & I_{r_1} \end{pmatrix}. \quad (27)$$

Introducing the matrix

$$W(z, \boldsymbol{\eta}) := \begin{pmatrix} 1 & \cdots & \kappa_{1,1}^{p_1+r_1-1} \\ \vdots & & \vdots \\ 1 & \cdots & \kappa_{1,p_1+r_1}^{p_1+r_1-1} \end{pmatrix}, \quad (28)$$

and using the previous expression of $\Pi^s(z, \boldsymbol{\eta})$, proving (27) is equivalent to proving the property

$$\lim_{z \rightarrow \infty} W(z, \boldsymbol{\eta})^{-1} \begin{pmatrix} I_{r_1} & 0 \\ 0 & 0 \end{pmatrix} W(z, \boldsymbol{\eta}) = \begin{pmatrix} I_{r_1} & 0 \\ 0 & 0 \end{pmatrix}. \quad (29)$$

We introduce the block decomposition of the matrix $W(z, \boldsymbol{\eta})$ in (28) and of its inverse $W(z, \boldsymbol{\eta})^{-1}$ (we forget temporarily to recall the $(z, \boldsymbol{\eta})$ dependence of all matrices to make expressions a little lighter):

$$W = \begin{pmatrix} V & V_{\sharp} \\ V_{\flat} & V_{\natural} \end{pmatrix}, \quad W^{-1} = \begin{pmatrix} \tilde{V} & \tilde{V}_{\sharp} \\ \tilde{V}_{\flat} & \tilde{V}_{\natural} \end{pmatrix},$$

where the upper left block has dimension $r_1 \times r_1$ and all other dimensions follow accordingly. The matrix on the left of (29), whose limit we wish to compute, is given by:

$$\begin{pmatrix} \tilde{V} V & \tilde{V} V_{\sharp} \\ \tilde{V}_b V & \tilde{V}_b V_{\sharp} \end{pmatrix} = \begin{pmatrix} I - \tilde{V}_{\sharp} V_b & \tilde{V} V_{\sharp} \\ -\tilde{V}_{\sharp} V_b & \tilde{V}_b V_{\sharp} \end{pmatrix},$$

where we have used the fact that $W^{-1}W$ is the identity. For computing the limit as z tends to infinity, we need some bounds on all matrices involved in the latter expression. Let us first examine the block V_{\sharp} . From the definition (28) of W and the labeling of eigenvalues of \mathbb{M} , the block V_{\sharp} involves powers at least equal to r_1 of stable eigenvalues of \mathbb{M} . The block V_b involves powers at most equal to $r_1 - 1$ of unstable eigenvalues of \mathbb{M} . Hence we have the bounds

$$\|V_{\sharp}\| = O(|z|^{-1}), \quad \|V_b\| = O(|z|^{(r_1-1)/p_1}),$$

where here $\|\cdot\|$ denotes any norm on (not necessarily square) complex matrices, for instance the maximum of the modulus of the entries. We now need to estimate the blocks of W^{-1} , which is made possible thanks to the explicit and somehow classical formula (the formula comes from the Lagrange polynomial interpolation theory):

$$(W^{-1})_{kj} = (-1)^k \left(\prod_{m \neq j} (\kappa_{1,m} - \kappa_{1,j}) \right)^{-1} \sum_{\substack{1 \leq m_1 < \dots < m_{p_1+r_1-k} \leq p_1+r_1, \\ m_1, \dots, m_{p_1+r_1-k} \neq j}} \kappa_{1,m_1} \cdots \kappa_{1,m_{p_1+r_1-k}}. \quad (30)$$

In this formula, the sum is understood as being equal to 1 if $k = p_1 + r_1$. Repeated and careful applications of (30) lead to the following bounds for the blocks of W^{-1} :

$$\|\tilde{V}\| = O(|z|^{1-1/r_1}), \quad \|\tilde{V}_{\sharp}\| = O(|z|^{-r_1/p_1-1/r_1}), \quad \|\tilde{V}_b\| = O(|z|^{1-1/p_1-1/r_1}), \quad \|\tilde{V}_{\sharp}\| = O(|z|^{-r_1/p_1}).$$

In particular, when combined with the bounds on V_{\sharp} and V_b , we can show that the four products of matrices $\tilde{V}_{\sharp} V_b$, $\tilde{V} V_{\sharp}$, $\tilde{V}_{\sharp} V_b$, and $\tilde{V}_b V_{\sharp}$ tend to zero, which completes the proof of (29).

Since we have (27), and the other limit

$$\lim_{z \rightarrow \infty} \Pi^u(z, \boldsymbol{\eta}) = \begin{pmatrix} I_{p_1} & 0 \\ 0 & 0 \end{pmatrix},$$

we get the limits of the stable and unstable subspaces as well. Namely, the limit $\mathbb{E}^s(\infty, \boldsymbol{\eta})$ of $\mathbb{E}^s(z, \boldsymbol{\eta})$ is the vector space spanned by the last r_1 vectors in the canonical basis of $\mathbb{C}^{p_1+r_1}$, and the limit $\mathbb{E}^u(\infty, \boldsymbol{\eta})$ of $\mathbb{E}^u(z, \boldsymbol{\eta})$ is the vector space spanned by the first p_1 vectors in the canonical basis of $\mathbb{C}^{p_1+r_1}$. With that definition, we obviously have the splitting (24) that persists up to $z = \infty$, as in the implicit case.

- Step 3. The transparent conditions in the physical variables.

In the previous step, we have seen that the projector $\Pi^s(z, \boldsymbol{\eta})$ onto the stable subspace of $\mathbb{M}(z, \boldsymbol{\eta})$ has a limit as z tends to infinity. For explicit schemes, this limit is independent of $\boldsymbol{\eta}$ and is given by (27). Consequently, Π^s extends as a function on $(\mathcal{U} \cup \{\infty\}) \times \mathbb{R}^{d-1}$ that depends holomorphically on z and analytically on $\boldsymbol{\eta}$ with the additional property of being 2π -periodic with respect to each coordinate of $\boldsymbol{\eta}$. We can therefore write the Laurent expansion of Π^s under the form

$$\forall (z, \boldsymbol{\eta}) \in (\mathcal{U} \cup \{\infty\}) \times \mathbb{R}^{d-1}, \quad \Pi^s(z, \boldsymbol{\eta}) = \sum_{n \geq 0} z^{-n} \Pi_n(\boldsymbol{\eta}),$$

where the convergence is normal on every compact subset. In particular, since each matrix Π_n depends analytically and in a periodic way on $\boldsymbol{\eta}$, the convergence of the Laurent series is normal on any set of the form $\{|z| \geq 1 + \delta\} \times \mathbb{R}^{d-1}$, $\delta > 0$, that is:

$$\forall \delta > 0, \quad \sum_{n \geq 0} \frac{1}{(1 + \delta)^n} \sup_{\boldsymbol{\eta} \in \mathbb{R}^{d-1}} \|\Pi_n(\boldsymbol{\eta})\| < +\infty. \quad (31)$$

We go back to the definition of the vector $U_{j_1}(\tau, \boldsymbol{\xi}')$ in (25), use the Laurent expansion of Π^s and the expression (19) of the Laplace-Fourier transform of u_{j_1} to get

$$\forall n \in \mathbb{N}, \quad \forall j_1 \leq 1, \quad \sum_{m=0}^n \Pi_{n-m}(\boldsymbol{\eta}) \begin{pmatrix} \widehat{u_{j_1+p_1-1}^m}(\boldsymbol{\xi}') \\ \vdots \\ \widehat{u_{j_1-r_1}^m}(\boldsymbol{\xi}') \end{pmatrix} = 0, \quad (32)$$

where we recall that here, the ‘hat’ denotes partial Fourier transform with respect to (x_2, \dots, x_d) , and $\boldsymbol{\eta}$ is a short notation for $(\xi_2 \Delta x_2, \dots, \xi_d \Delta x_d)$. The sequence of operators $(\Pi_n)_{n \in \mathbb{N}}$ is then defined as the following sequence of Fourier multipliers: for any $n \in \mathbb{N}$ and any sequence $v \in \ell^2(\mathbb{Z}^{d-1}; \mathbb{C}^{p_1+r_1})$, we identify the sequence v and the corresponding step function

$$v(x') := v_{j'}, \quad \forall x' \in \prod_{k=2}^d [j_k \Delta x_k; (j_k + 1) \Delta x_k).$$

In particular, the Fourier transform of the sequence v means the Fourier transform of the corresponding step function. Then $\Pi_n v$ is defined as the sequence⁸ whose Fourier transform is given by

$$\boldsymbol{\xi}' \in \mathbb{R}^{d-1} \mapsto \Pi_n(\boldsymbol{\eta}) \widehat{v}(\boldsymbol{\xi}'). \quad (33)$$

With this definition for the Π_n ’s, applying the inverse Fourier transform to (32) gives (11). To complete the proof of Theorem 1, it only remains to show that the Π_n ’s satisfy the growth condition and the algebraic constraints stated in Theorem 1. The growth condition is a direct consequence of the bound (31) on the symbols (Π_n) of the Fourier multipliers (Π_n) . The algebraic constraints follow by using the fact that Π^s is a projector (up to now we have only used that E^u is the kernel of Π^s). Expanding the equality $\Pi^s(z, \boldsymbol{\eta})^2 = \Pi^s(z, \boldsymbol{\eta})$ in Laurent series with respect to z , we get the algebraic constraints for the symbols:

$$\forall n \in \mathbb{N}, \quad \forall \boldsymbol{\eta} \in \mathbb{R}^{d-1}, \quad \Pi_n(\boldsymbol{\eta}) = \sum_{m=0}^n \Pi_m(\boldsymbol{\eta}) \Pi_{n-m}(\boldsymbol{\eta}).$$

The relations satisfied by the Π_n ’s follow immediately. This completes the proof of Theorem 1.

2.2 Alternative formulations of transparent boundary conditions

In this paragraph, we explain why the formulation of the transparent conditions encoded in the operators Π_n in (11) is not unique and what other choices, that may be more suitable from a practical point of view, can be made.

⁸It is indeed a rather standard result that Fourier multipliers associated with periodic symbols map the set of L^2 step functions into itself. We can therefore equivalently view the operator Π_n as acting on $\ell^2(\mathbb{Z}^{d-1})$ with values in $\ell^2(\mathbb{Z}^{d-1})$ rather than acting on the set of L^2 step functions with values in $L^2(\mathbb{R}^{d-1})$.

2.2.1 Alternative formulation with projectors

In the Laplace-Fourier variables, the recurrence relation (23) implies that the vector $U_1(\tau, \xi')$ belongs to the unstable subspace $\mathbb{E}^u(z, \boldsymbol{\eta})$ of $\mathbb{M}(z, \boldsymbol{\eta})$. Using the decomposition (24), we have equivalently formulated this property in writing (25). The arbitrariness here lies in the choice of the supplementary vector space $\mathbb{E}^s(z, \boldsymbol{\eta})$. More precisely, let us assume that one can choose a vector space $\tilde{\mathbb{E}}^s(z, \boldsymbol{\eta})$ of dimension r_1 in $\mathbb{C}^{p_1+r_1}$, depending holomorphically on $z \in \mathcal{U} \cup \{\infty\}$, analytically on $\boldsymbol{\eta} \in \mathbb{R}^{d-1}$ and 2π -periodically with respect to each coordinate of $\boldsymbol{\eta}$, and such that one has the decomposition⁹:

$$\forall (z, \boldsymbol{\eta}) \in (\mathcal{U} \cup \{\infty\}) \times \mathbb{R}^{d-1}, \quad \mathbb{C}^{p_1+r_1} = \tilde{\mathbb{E}}^s(z, \boldsymbol{\eta}) \oplus \mathbb{E}^u(z, \boldsymbol{\eta}). \quad (34)$$

We then let $\tilde{\Pi}^{s,u}(z, \boldsymbol{\eta})$ denote the corresponding projectors¹⁰. Then one can equivalently rewrite (25) as

$$\forall j_1 \leq 1, \quad \tilde{\Pi}^s(z, \boldsymbol{\eta}) U_{j_1}(\tau, \xi') = 0,$$

and since we already know that $\tilde{\Pi}^s$ is holomorphic in z on $\mathcal{U} \cup \{\infty\}$, the Laurent series of $\tilde{\Pi}^s$ only involves nonpositive powers of z . We can therefore reproduce the same arguments as in Step 3 of the proof of Theorem 1, which gives rise in the end to a set of relations

$$\forall n \in \mathbb{N}, \quad \forall j_1 \leq 0, \quad \sum_{m=0}^n \tilde{\Pi}_{n-m} \begin{pmatrix} u_{(j_1+p_1, \cdot)}^m \\ \vdots \\ u_{(j_1+1-r_1, \cdot)}^m \end{pmatrix} = 0,$$

with a definition of the operators $\tilde{\Pi}_{n-m}$ that is entirely analogous to the one of the Π_n 's. What really matters is not the sequence of operators (Π_n) in (11) but rather the kernel and the range of Π_0 for instance. This is the reason why there is some possible freedom in the formulation of (11).

Of course, the supplementary vector space $\mathbb{E}^s(z, \boldsymbol{\eta})$ in (24) is a rather natural choice, at least because when the eigenvalues of $\mathbb{M}(z, \boldsymbol{\eta})$ are simple, the projectors $\Pi^{s,u}(z, \boldsymbol{\eta})$ are given in terms of a Vandermonde matrix and its inverse for which explicit expressions are available¹¹. Therefore there does not seem to be much simplification in choosing another supplementary vector space to \mathbb{E}^u , though one should keep in mind that it is a possibility.

2.2.2 Alternative formulation with linear forms

The other formulation that we propose seems to be much more used in practice, see e.g. [EA01, AES03, ZE06, Ehr10, DZ06, ZWH08, BELV16]. It is specifically recommended in the case $r_1 = 1$ for then $\mathbb{E}^u(z, \boldsymbol{\eta})$ is a hyperplane in $\mathbb{C}^{p_1+r_1}$ which one can consider as the kernel of some linear form.

In full generality, we know that \mathbb{E}^u defines a holomorphic/analytic-periodic vector bundle over $(\mathcal{U} \cup \{\infty\}) \times \mathbb{R}^{d-1}$. By holomorphic/analytic-periodic, it should be clear by now that we mean holomorphic with respect to $z \in \mathcal{U} \cup \{\infty\}$, analytic with respect to $\boldsymbol{\eta} \in \mathbb{R}^{d-1}$ and 2π -periodic with respect to each coordinate of $\boldsymbol{\eta}$. In the same way, \mathbb{E}^s defines a holomorphic/analytic-periodic vector bundle over $(\mathcal{U} \cup \{\infty\}) \times \mathbb{R}^{d-1}$, and the fiber $\mathbb{E}^s(z, \boldsymbol{\eta})$, resp. $\mathbb{E}^u(z, \boldsymbol{\eta})$, coincides with the range of the projector $\Pi^s(z, \boldsymbol{\eta})$, resp. $\Pi^u(z, \boldsymbol{\eta})$. By the transformation $z \rightarrow 1/z$, the set $\mathcal{U} \cup \{\infty\}$ is mapped biholomorphically onto the unit disk \mathbb{D} , which is simply connected. By following the argument in [Kat95, Chapter 2.4], we can thus determine

⁹We recall here that \mathbb{E}^u has a limit when z tends to infinity, which we have denoted $\mathbb{E}^u(\infty, \boldsymbol{\eta})$.

¹⁰Of course, $\tilde{\Pi}^u(z, \boldsymbol{\eta})$ does not coincide with $\Pi^u(z, \boldsymbol{\eta})$, even though one of the vector spaces in (34) is the same as in (24).

¹¹There are also explicit expressions when some eigenvalues are not simple but the algebra gets more involved.

for all $\boldsymbol{\eta} \in \mathbb{R}^{d-1}$ a basis $e_1^s(z, \boldsymbol{\eta}), \dots, e_{r_1}^s(z, \boldsymbol{\eta})$ of $\mathbb{E}^s(z, \boldsymbol{\eta})$, resp. a basis $e_1^u(z, \boldsymbol{\eta}), \dots, e_{p_1}^u(z, \boldsymbol{\eta})$ of $\mathbb{E}^u(z, \boldsymbol{\eta})$, that depends holomorphically on $z \in \mathcal{U} \cup \{\infty\}$. (We do not concentrate at first on the dependence with respect to $\boldsymbol{\eta}$ and consider $\boldsymbol{\eta}$ as a fixed parameter for now.)

The inverse of the matrix

$$(e_1^s(z, \boldsymbol{\eta}) \cdots e_{r_1}^s(z, \boldsymbol{\eta}) e_1^u(z, \boldsymbol{\eta}) \cdots e_{p_1}^u(z, \boldsymbol{\eta})) \in \mathcal{M}_{p_1+r_1}(\mathbb{C}),$$

provides with r_1 row vectors $L_1(z, \boldsymbol{\eta}), \dots, L_{r_1}(z, \boldsymbol{\eta}) \in \mathcal{M}_{1, p_1+r_1}(\mathbb{C})$ that depend holomorphically on $z \in \mathcal{U} \cup \{\infty\}$ and such that the unstable subspace \mathbb{E}^u reads

$$\mathbb{E}^u(z, \boldsymbol{\eta}) = \{X \in \mathbb{C}^{p_1+r_1} / L_1(z, \boldsymbol{\eta}) X = \cdots = L_{r_1}(z, \boldsymbol{\eta}) X = 0\}.$$

In other words, for all $\boldsymbol{\eta} \in \mathbb{R}^{d-1}$, we have constructed a matrix $L(z, \boldsymbol{\eta}) \in \mathcal{M}_{r_1, p_1+r_1}(\mathbb{C})$ of *full rank* r_1 , that depends holomorphically on z on $\mathcal{U} \cup \{\infty\}$, and whose kernel is $\mathbb{E}^u(z, \boldsymbol{\eta})$. The relations (25) in the proof of Theorem 1 can be equivalently recast as

$$\forall j_1 \leq 1, \quad L(z, \boldsymbol{\eta}) U_{j_1}(\tau, \boldsymbol{\xi}') = 0,$$

with the major gain, from an algebraic point of view, that L has full rank so that any perturbation (meaning a nonzero vector in \mathbb{C}^{r_1} on the right hand side is now admissible in view of a stability analysis, which bypasses the algebraic constraints of Lemma 2). We can expand L in Laurent series

$$L(z, \boldsymbol{\eta}) = \sum_{n \geq 0} \frac{1}{z^n} L_n(\boldsymbol{\eta}),$$

apply the inverse Laplace transform and rewrite the above relations as

$$\forall n \in \mathbb{N}, \quad \forall j_1 \leq 1, \quad \sum_{m=0}^n L_{n-m}(\boldsymbol{\eta}) \begin{pmatrix} \widehat{u_{j_1+p_1-1}^m}(\boldsymbol{\xi}') \\ \vdots \\ \widehat{u_{j_1-r_1}^m}(\boldsymbol{\xi}') \end{pmatrix} = 0,$$

which is the analogue of (32).

Assume now that the above construction of the matrix L can be performed in such a way that L depends analytically and 2π -periodically on $\boldsymbol{\eta}$. Analyticity can be obtained by using the same procedure as in [BGS07, Chapter 4.6] (which corresponds to applying the method of [Kat95] coordinate by coordinate), but periodicity seems far from obvious if one constructs the above basis by ODE arguments as in [Kat95]. Assume nevertheless that periodicity can be achieved (for instance, as in the example below, because one has explicit expressions). Then one has all the desirable properties for applying inverse Fourier transform in the previous relations and write the relations in the original physical variables under the form

$$\forall n \in \mathbb{N}, \quad \forall j_1 \leq 0, \quad \sum_{m=0}^n \mathbf{L}_{n-m} \begin{pmatrix} u_{(j_1+p_1, \cdot)}^m \\ \vdots \\ u_{(j_1+1-r_1, \cdot)}^m \end{pmatrix} = 0,$$

with suitable Fourier multipliers \mathbf{L}_n on $\ell^2(\mathbb{Z}^{d-1})$. The nice feature here is that one no longer needs to care about compatibility conditions for the boundary forcing terms (if present). There may of course still remain some indeterminacy in the construction of L due again to the choice of a supplementary vector space to \mathbb{E}^u and to the choice of a basis.

In the case $d = 1$, the argument in [Kat95] allows us to construct the matrix $L(z)$ as above, with L holomorphic on $\mathcal{U} \cup \{\infty\}$. Such a construction is the one that is used for instance in [EA01, BELV16, BMGN16]. When d is larger than 1, the main difficulty is to construct $L(z, \boldsymbol{\eta})$ with the property that L is 2π -periodic with respect to each coordinate of $\boldsymbol{\eta}$. If we go through the arguments above, we could construct such an L provided that we have a basis of \mathbb{E}^s and a basis of \mathbb{E}^u that depend holomorphically/analytically/periodically on $(z, \boldsymbol{\eta})$ in $(\mathcal{U} \cup \{\infty\}) \times \mathbb{R}^{d-1}$. Constructing such bases is far from obvious, because it amounts to showing that the vector bundles \mathbb{E}^s and \mathbb{E}^u are trivial, but here, because of the periodicity in $\boldsymbol{\eta}$, these vector bundles are considered over $(\mathcal{U} \cup \{\infty\}) \times (\mathbb{S}^1)^{d-1}$, which is *not contractible* (unless $d = 1$). However, there are examples for which one can show from an explicit expression of $\mathbb{E}^s(z, \boldsymbol{\eta})$ and $\mathbb{E}^u(z, \boldsymbol{\eta})$ that the bundles over $(\mathcal{U} \cup \{\infty\}) \times (\mathbb{S}^1)^{d-1}$ are indeed trivial, and this allows then to construct the matrix $L(z, \boldsymbol{\eta})$ even for problems with $d \geq 2$.

2.2.3 The simplest and most frequent example

Let us consider for simplicity the case $d = 1$ so that the previous technical difficulties encountered because of the tangential Fourier variables $\boldsymbol{\eta}$ disappear. Let us assume furthermore $p_1 = r_1 = 1$. In other words, the original numerical scheme (3) takes the form:

$$\begin{cases} \sum_{\sigma=0}^{s+1} a_{-1,\sigma} u_{j-1}^{n+\sigma} + a_{0,\sigma} u_j^{n+\sigma} + a_{1,\sigma} u_{j+1}^{n+\sigma} = 0, & n \geq 0, j \in \mathbb{Z}, \\ (u^0, \dots, u^s) = (f^0, \dots, f^s) \in \ell^2(\mathbb{Z})^{s+1}. \end{cases} \quad (35)$$

Of course, the coefficients $a_{-1,\sigma}, a_{0,\sigma}, a_{1,\sigma}$ in (35) may depend on Δt and/or Δx .

Let us emphasize how the various assumptions on (3) translate in the particular case (35). Assumptions 1 and 5 mean:

For explicit schemes: $a_{-1,s+1} = a_{1,s+1} = 0$, $a_{0,s+1} = 1$, and $a_{-1,s} a_{1,s} \neq 0$ (this last condition may restrict the possible values of the discretization parameters $\Delta t, \Delta x$, see the example of the Lax-Wendroff scheme in Section 5).

For implicit schemes: $a_{-1,s+1} a_{1,s+1} \neq 0$, and the polynomial (in κ)

$$a_{-1,s+1} + a_{0,s+1} \kappa + a_{1,s+1} \kappa^2,$$

has one root in \mathbb{D} (necessarily not zero) and one root in \mathcal{U} . If the coefficients are real, this means that one root belongs to $(-1, 1)$ and one root belongs to $(-\infty, -1) \cup (1, +\infty)$ for the roots cannot be complex conjugate.

Assumption 3 means that both polynomials

$$\sum_{\sigma=0}^{s+1} a_{-1,\sigma} z^\sigma, \quad \sum_{\sigma=0}^{s+1} a_{1,\sigma} z^\sigma,$$

have no root in \mathcal{U} (Assumption 4 means that they have no root in $\overline{\mathcal{U}}$). For explicit schemes, these polynomials have degree s while they have degree $s + 1$ for implicit schemes. In particular, Assumption 4 is automatically satisfied if the scheme is explicit and $s = 0$ for then the previous two polynomials are constant and nonzero.

Let us now turn to our derivation of the transparent boundary conditions. The matrix \mathbb{M} of interest is defined in (22). In one space dimension with $p = r = 1$, its expression reduces to

$$\mathbb{M}(z) := \begin{pmatrix} -\frac{a_0(z)}{a_1(z)} & -\frac{a_{-1}(z)}{a_1(z)} \\ 1 & 0 \end{pmatrix} \in \mathcal{M}_2(\mathbb{C}), \quad a_\ell(z) := \sum_{\sigma=0}^{s+1} a_{\ell,\sigma} z^\sigma. \quad (36)$$

We know from the proof of Theorem 1 that $\mathbb{M}(z)$ has one eigenvalue $\kappa_s(z)$ in \mathbb{D} and one eigenvalue $\kappa_u(z)$ in \mathcal{U} , both counted with multiplicity and therefore simple. Since the eigenvalues do not cross for $z \in \mathcal{U}$, they both depend holomorphically on $z \in \mathcal{U}$. The stable eigenvalue $\kappa_s(z)$ has a limit when z tends to infinity, while the unstable eigenvalue $\kappa_u(z)$ has a limit when z tends to infinity only when the scheme is implicit (for an explicit scheme, $\kappa_u(z)$ tends to infinity when z tends to infinity). The stable and unstable subspaces read

$$\mathbb{E}^s(z) = \text{Span} \begin{pmatrix} \kappa_s(z) \\ 1 \end{pmatrix}, \quad \mathbb{E}^u(z) = \text{Span} \begin{pmatrix} 1 \\ \kappa_u(z)^{-1} \end{pmatrix},$$

where the choice for parametrizing \mathbb{E}^u has been made in such a way that the generating vector has a limit when z tends to infinity (even when the scheme is explicit).

Let us now discuss two possible ways of writing the transparent boundary conditions. We first follow the approach based on the spectral projectors as in the proof of Theorem 1. The decomposition:

$$\forall z \in \mathcal{U} \cup \{\infty\}, \quad \mathbb{C}^2 = \mathbb{E}^s(z) \oplus \mathbb{E}^u(z),$$

is endowed with two projectors $\Pi^{s,u}(z)$, whose explicit expression is given by

$$\Pi^s(z) = \frac{1}{\kappa_s(z) - \kappa_u(z)} \begin{pmatrix} \kappa_s(z) & -\kappa_s(z) \kappa_u(z) \\ 1 & -\kappa_u(z) \end{pmatrix}, \quad \Pi^u(z) = I - \Pi^s(z).$$

Observe in particular that for explicit schemes, knowing $\kappa_s \rightarrow 0$ and $\kappa_u \rightarrow \infty$ as z tends to ∞ , we recover the asymptotic behavior

$$\lim_{z \rightarrow \infty} \Pi^s(z) = \begin{pmatrix} 0 & 0 \\ 0 & 1 \end{pmatrix}, \quad \lim_{z \rightarrow \infty} \Pi^u(z) = \begin{pmatrix} 1 & 0 \\ 0 & 0 \end{pmatrix}.$$

We can either write the vector space $\mathbb{E}^u(z)$ as the kernel of the matrix $\Pi^s(z)$ or as the kernel of the linear form

$$x = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \in \mathbb{C}^2 \mapsto \frac{1}{\kappa_u(z)} x_1 - x_2.$$

Again, it is more convenient to use $\kappa_u(z)^{-1}$ instead of $\kappa_u(z)$ because in our framework, $\kappa_u(z)^{-1}$ always has a limit as z tends to infinity.

The Laurent series of both κ_s and κ_u^{-1} can be obtained by starting from the relations

$$a_{-1}(z) + a_0(z) \kappa_s(z) + a_1(z) \kappa_s(z)^2 = 0, \quad a_{-1}(z) \kappa_u(z)^{-2} + a_0(z) \kappa_u(z)^{-1} + a_1(z) = 0,$$

and by identifying inductively the coefficients in the series

$$\kappa_s(z) = \sum_{n \geq 0} \frac{\kappa_{s,n}}{z^n}, \quad \kappa_u(z)^{-1} = \sum_{n \geq 0} \frac{\tilde{\kappa}_{u,n}}{z^n}.$$

For the numerical scheme considered in [EA01], one can get an explicit expression of the coefficients in terms of the Legendre polynomials. This expression does not come from the general method we propose. For explicit schemes, both $\kappa_{s,0}$ and $\tilde{\kappa}_{u,0}$ are zero, and both $\kappa_{s,1}$ and $\tilde{\kappa}_{u,1}$ are nonzero. One can also determine the Laurent series expansion of κ_u by writing (the coefficient $\kappa_{u,1}$ being nonzero for explicit schemes):

$$\kappa_u(z) = \kappa_{u,1} z + \sum_{n \geq 0} \frac{\kappa_{u,n}}{z^n},$$

by plugging this series into the equation

$$a_{-1}(z) + a_0(z) \kappa_u(z) + a_1(z) \kappa_u(z)^2 = 0,$$

and by identifying inductively the coefficients $\kappa_{u,n}$. At least theoretically speaking, this gives access to the Laurent series expansion of any useful quantity involving κ_s and κ_u . In particular, we have access to the Laurent series of Π^s , which gives rise in the original physical variables to the relations (recall here that we consider a one-dimensional problem so the Fourier multipliers Π_n reduce in fact to 2×2 matrices):

$$\forall n \in \mathbb{N}, \quad \forall j_1 \leq 0, \quad \sum_{m=0}^n \Pi_{n-m} \begin{pmatrix} u_{j_1+1}^m \\ u_{j_1}^m \end{pmatrix} = 0.$$

Viewing the vector space $\mathbb{E}^u(z)$ as the kernel of a linear form, we may use the Laurent expansion of κ_u^{-1} and rewrite equivalently the latter relations as

$$\forall n \in \mathbb{N}, \quad \forall j_1 \leq 0, \quad u_{j_1}^n = \sum_{m=0}^n \tilde{\kappa}_{u,n-m} u_{j_1+1}^m. \quad (37)$$

Using the formulation (37), the new (though equivalent) way of writing the transparent boundary conditions in (12) becomes

$$\begin{cases} \sum_{\sigma=0}^{s+1} a_{-1,\sigma} u_{j-1}^{n+\sigma} + a_{0,\sigma} u_j^{n+\sigma} + a_{1,\sigma} u_{j+1}^{n+\sigma} = \Delta t F_j^{n+s+1}, & n \geq 0, \quad j \geq 1, \\ u_0^{n+s+1} - \tilde{\kappa}_{u,0} u_1^{n+s+1} = \sum_{m=0}^{n+s} \tilde{\kappa}_{u,n+s+1-m} u_1^m + g^{n+s+1}, & n \geq 0, \\ (u_j^0, \dots, u_j^s) = (f_j^0, \dots, f_j^s), & j \geq 0. \end{cases} \quad (38)$$

For explicit schemes, $\tilde{\kappa}_{u,0}$ is zero so the boundary condition in (38) takes the form of a non-homogeneous Dirichlet boundary condition, whose source term is computed thanks to the trace at $j = 1$ of the solution at earlier times.

3 Solvability of the scheme with transparent boundary conditions

From now on, we analyze the numerical scheme (12). We first prove Lemma 2 which characterizes the sequence of source terms for which one can construct a solution to the sequence of ‘convolution’ equations (13). Lemma 2 will provide with necessary and sufficient compatibility conditions for solving (12).

Proof of Lemma 2. The whole proof is based on the mere fact that P_0 is a projector. Hence a vector y belongs to the range of P_0 if and only if $P_0 y = y$.

Let us first assume that the sequence (y_n) is such that there exists a solution (x_n) to (13). For $n = 0$, there holds $y_0 \in \text{Im } P_0$ and therefore (14) holds for $n = 0$. Let us assume by induction that (14) holds up to some integer n . Then we can write

$$y_{n+1} - \sum_{m=0}^n P_{n+1-m} x_m = P_0 x_{n+1} \in \text{Im } P_0,$$

and consequently

$$\begin{aligned} y_{n+1} - P_0 y_{n+1} &= \sum_{m=0}^n P_{n+1-m} x_m - \sum_{m=0}^n P_0 P_{n+1-m} x_m \\ &= (P_{n+1} - P_0 P_{n+1}) x_0 + \sum_{m=1}^n (P_{n+1-m} - P_0 P_{n+1-m}) x_m. \end{aligned} \quad (39)$$

We now use the equations satisfied by the P_n 's and write

$$\begin{aligned} (P_{n+1} - P_0 P_{n+1}) x_0 &= P_{n+1} P_0 x_0 + \sum_{p=1}^n P_p P_{n+1-p} x_0 \\ &= P_{n+1} y_0 + \sum_{p=1}^n P_p P_{n+1-p} x_0. \end{aligned}$$

We can therefore simplify (39) and obtain

$$y_{n+1} - P_0 y_{n+1} - P_{n+1} y_0 = \sum_{m=1}^n (P_{n+1-m} - P_0 P_{n+1-m}) x_m + \sum_{p=1}^n P_p P_{n+1-p} x_0. \quad (40)$$

If $n = 0$, the work is over since the sums on the right hand side of (40) vanish and we have obtained (14) for $n + 1 = 1$. Assuming $n \geq 1$ from now on, we substitute $P_{n+1-p} x_0$ in (40), $p = 1, \dots, n$, for

$$y_{n+1-p} - \sum_{k=0}^{n-p} P_{n+1-p-k} x_k,$$

by using (13). Using the algebraic relations satisfied by the P_k 's and some manipulations on the indices, we end up with

$$y_{n+1} - P_0 y_{n+1} - P_{n+1} y_0 = \sum_{p=1}^n P_p y_{n+1-p},$$

which is nothing but (14) for $n + 1$.

Let us now assume that the sequence (y_n) satisfies the compatibility conditions (14). Then one immediately sees that the sequence (x_n) defined by $x_n := y_n$ for all n satisfies (13). The proof of Lemma 2 is thus complete. \square

It remains to use Lemma 2 for proving the solvability result of Proposition 1.

Proof of Proposition 1. We first show that if there exists a solution to (12), then it is necessarily unique. Of course, by linearity, this amounts to showing that the only solution to

$$\begin{cases} \sum_{\sigma=0}^{s+1} Q_{\sigma} u_j^{n+\sigma} = 0, & n \geq 0, \quad j_1 \geq 1, \\ \sum_{m=0}^{n+s+1} \mathbf{\Pi}_{n+s+1-m} \begin{pmatrix} u_{(p_1, \cdot)}^m \\ \vdots \\ u_{(1-r_1, \cdot)}^m \end{pmatrix} = 0, & n \geq 0, \\ (u_j^0, \dots, u_j^s) = (0, \dots, 0), & j_1 \geq 1 - r_1, \end{cases} \quad (41)$$

is zero. We prove this property by induction on the time level. Let us assume that up to some level $n \geq 0$, there holds $u^0 = \dots = u^{n+s} = 0$ (this property clearly holds for $n = 0$ due to the initial data in (41)). Then $u^{n+s+1} \in \ell^2$ is a solution to

$$\begin{cases} Q_{s+1} u_j^{n+s+1} = 0, & j_1 \geq 1, \\ \mathbf{\Pi}_0 \begin{pmatrix} u_{(p_1, \cdot)}^{n+s+1} \\ \vdots \\ u_{(1-r_1, \cdot)}^{n+s+1} \end{pmatrix} = 0. \end{cases}$$

We apply a partial Fourier transform with respect to the tangential variables $j' \in \mathbb{Z}^{d-1}$, which yields, because $\mathbf{\Pi}_0$ is a Fourier multiplier, see (33):

$$\begin{cases} Q_{s+1}^{\sharp}(\boldsymbol{\eta}) \widehat{u_{j_1}^{n+s+1}}(\boldsymbol{\xi}') = 0, & j_1 \geq 1, \\ \mathbf{\Pi}_0(\boldsymbol{\eta}) \begin{pmatrix} \widehat{u_{p_1}^{n+s+1}}(\boldsymbol{\xi}') \\ \vdots \\ \widehat{u_{1-r_1}^{n+s+1}}(\boldsymbol{\xi}') \end{pmatrix} = 0, \end{cases}$$

where the ‘one-dimensional’ finite difference operator Q_{s+1}^{\sharp} is defined in (20).

In the explicit case (case (i) in Assumption 5), we compute $Q_{s+1}^{\sharp}(\boldsymbol{\eta}) = I$, and we have seen in the proof of Theorem 1 that the projector $\mathbf{\Pi}_0(\boldsymbol{\eta})$ reduces to

$$\mathbf{\Pi}_0(\boldsymbol{\eta}) = \begin{pmatrix} 0 & 0 \\ 0 & I_{r_1} \end{pmatrix}.$$

Hence we get $\widehat{u_{j_1}^{n+s+1}}(\boldsymbol{\xi}') = 0$ for all $j_1 \geq 1 - r_1$, and applying the inverse Fourier transform, we obtain $u^{n+s+1} = 0$. Uniqueness follows.

In the implicit case (case (ii) in Assumption 5), we know from the index condition (5) that a sequence $(w_{j_1})_{j_1 \geq 1-r_1}$ solution to the recurrence relation

$$\forall j_1 \geq 1, \quad Q_{s+1}^{\sharp}(\boldsymbol{\eta}) w_{j_1} = 0, \quad (42)$$

belongs to ℓ^2 if and only if its ‘initial condition’ $(w_{p_1}, \dots, w_{1-r_1})^T$ belongs to the stable subspace $\mathbb{E}^s(\infty, \boldsymbol{\eta})$ of the matrix $\mathbb{M}(\infty, \boldsymbol{\eta})$ whose expression is given in (26). The matrix $\mathbb{M}(\infty, \boldsymbol{\eta})$ is of course the companion

matrix that arises when one rewrites the recurrence (42) as a first order recurrence relation for the augmented vector $(w_{j_1+p_1-1}, \dots, w_{j_1-r_1})^T$. Moreover, we know from the proof of Theorem 1 that the projector $\Pi_0(\boldsymbol{\eta})$, which is the limit of $\Pi^s(z, \boldsymbol{\eta})$ as z tends to infinity, is precisely the projection on $\mathbb{E}^s(\infty, \boldsymbol{\eta})$ with kernel $\mathbb{E}^u(\infty, \boldsymbol{\eta})$. Since the sequence u^{n+s+1} belongs to $\ell^2([1-r_1, +\infty) \times \mathbb{Z}^{d-1})$, the partial Fourier transform $(\widehat{u_{j_1}^{n+s+1}}(\boldsymbol{\xi}'))_{j_1 \geq 1-r_1}$ belongs to ℓ^2 for almost every $\boldsymbol{\xi}'$. We therefore get $(\widehat{u_{p_1}^{n+s+1}}(\boldsymbol{\xi}'), \dots, \widehat{u_{1-r_1}^{n+s+1}}(\boldsymbol{\xi}')) = 0$ for almost every $\boldsymbol{\xi}'$, which yields $u^{n+s+1} = 0$ by inverse Fourier transform. Uniqueness of a solution to (12) follows by induction on n .

It remains to show that there exists a solution to (12), provided that the necessary compatibility conditions described in Proposition 1 are satisfied. We thus assume that the source terms g^n in (12) satisfy (15) (for $n \geq s+1$, g^n refers to the boundary forcing term in (12), and for $0 \leq n \leq s$, g^n refers to the sequence constructed from the initial data as in the statement of Proposition 1). Let us assume that up to some time level $n+s$, with $n \geq 0$, we have constructed the sequences u^0, \dots, u^{n+s} , with

$$\begin{aligned} (u^0, \dots, u^s) &= (f^0, \dots, f^s), \quad (\text{initial data}), \\ \forall n' = 1, \dots, n, \quad \sum_{k=0}^{n'+s} \mathbf{\Pi}_{n'+s-k} \begin{pmatrix} u_{(p_1, \cdot)}^m \\ \vdots \\ u_{(1-r_1, \cdot)}^m \end{pmatrix} &= g^{n'+s}, \quad (\text{boundary conditions}), \\ \forall n' = 1, \dots, n, \quad \forall j_1 \geq 1, \quad \sum_{\sigma=0}^{s+1} Q_\sigma u_j^{n'+\sigma-1} &= \Delta t F_j^{n'+s}, \quad (\text{numerical scheme}). \end{aligned}$$

We wish to construct a sequence u^{n+s+1} solution to

$$\begin{cases} Q_{s+1} u_j^{n+s+1} = \Delta t F_j^{n+s+1} - \sum_{\sigma=0}^s Q_\sigma u_j^{n+\sigma}, & j_1 \geq 1, \\ \mathbf{\Pi}_0 \begin{pmatrix} u_{(p_1, \cdot)}^{n+s+1} \\ \vdots \\ u_{(1-r_1, \cdot)}^{n+s+1} \end{pmatrix} = g^{n+s+1} - \sum_{m=0}^{n+s} \mathbf{\Pi}_{n+s+1-m} \begin{pmatrix} u_{(p_1, \cdot)}^m \\ \vdots \\ u_{(1-r_1, \cdot)}^m \end{pmatrix}. \end{cases} \quad (43)$$

Let us first observe that, using Lemma 2, we already know that the sequence

$$g^{n+s+1} - \sum_{m=0}^{n+s} \mathbf{\Pi}_{n+s+1-m} \begin{pmatrix} u_{(p_1, \cdot)}^m \\ \vdots \\ u_{(1-r_1, \cdot)}^m \end{pmatrix},$$

belongs to the range of the projector $\mathbf{\Pi}_0$ (one can for instance reproduce similar calculations as in the proof of Lemma 2 and prove that this sequence belongs to the kernel of $\mathbf{\Pi}_0 - I$). We can therefore look for the solution u^{n+s+1} to (43) under the form $u^{n+s+1} = v + w$, with

$$\begin{pmatrix} v_{(p_1, \cdot)} \\ \vdots \\ v_{(1-r_1, \cdot)} \end{pmatrix} := g^{n+s+1} - \sum_{m=0}^{n+s} \mathbf{\Pi}_{n+s+1-m} \begin{pmatrix} u_{(p_1, \cdot)}^m \\ \vdots \\ u_{(1-r_1, \cdot)}^m \end{pmatrix},$$

and for instance $v_{(j_1, \cdot)} := 0$ if $j_1 > p_1$. We are then reduced to showing that, for some sequence $(\tilde{F}_j) \in \ell^2$ whose expression is not useful, there exists a sequence $w \in \ell^2$ solution to

$$\begin{cases} Q_{s+1} w_j = \tilde{F}_j, & j_1 \geq 1, \\ \Pi_0 \begin{pmatrix} w_{(p_1, \cdot)} \\ \vdots \\ w_{(1-r_1, \cdot)} \end{pmatrix} = 0. \end{cases} \quad (44)$$

For simplicity, we drop the tilde on the source term in (44). We construct the partial Fourier transform of the solution w to (44) rather than w itself. Namely, we are going to construct a solution to

$$\begin{cases} Q_{s+1}^\#(\boldsymbol{\eta}) \hat{w}_{j_1}(\boldsymbol{\xi}') = \hat{F}_{j_1}(\boldsymbol{\xi}'), & j_1 \geq 1, \\ \Pi_0(\boldsymbol{\eta}) \begin{pmatrix} \hat{w}_{p_1}(\boldsymbol{\xi}') \\ \vdots \\ \hat{w}_{1-r_1}(\boldsymbol{\xi}') \end{pmatrix} = 0. \end{cases} \quad (45)$$

Let us first deal with the explicit case. In that case, one simply has to set

$$\hat{w}_{j_1}(\boldsymbol{\xi}') := \begin{cases} 0, & \text{if } j_1 \leq 0, \\ \hat{F}_{j_1}(\boldsymbol{\xi}'), & \text{if } j_1 \geq 1, \end{cases}$$

apply inverse Fourier transform and obtain a solution $w \in \ell^2$ to (44). We therefore focus on the implicit case. With the vector $\mathcal{W}_{j_1}(\boldsymbol{\xi}') := (\hat{w}_{j_1+p_1-1}(\boldsymbol{\xi}'), \dots, \hat{w}_{j_1-r_1}(\boldsymbol{\xi}'))^T$, the system (45) is equivalently rewritten as

$$\begin{cases} \mathcal{W}_{j_1+1}(\boldsymbol{\xi}') - \mathbb{M}(\infty, \boldsymbol{\eta}) \mathcal{W}_{j_1}(\boldsymbol{\xi}') = a_{\infty, p_1}(\boldsymbol{\eta})^{-1} (\hat{F}_{j_1}(\boldsymbol{\xi}'), 0, \dots, 0)^T, & j_1 \geq 1, \\ \mathcal{W}_1(\boldsymbol{\xi}') \in \mathbb{E}^u(\infty, \boldsymbol{\eta}). \end{cases} \quad (46)$$

The unique ℓ^2 solution to this problem is written explicitly by using the spectral projectors $\Pi_0(\boldsymbol{\eta})$ and $I - \Pi_0(\boldsymbol{\eta})$. We obtain the formula

$$\begin{aligned} \mathcal{W}_{j_1}(\boldsymbol{\xi}') &= a_{\infty, p_1}(\boldsymbol{\eta})^{-1} \sum_{k=1}^{j_1-1} \left(\mathbb{M}(\infty, \boldsymbol{\eta}) \Pi_0(\boldsymbol{\eta}) \right)^{j_1-1-k} (\hat{F}_{j_1}(\boldsymbol{\xi}'), 0, \dots, 0)^T \\ &\quad - a_{\infty, p_1}(\boldsymbol{\eta})^{-1} \sum_{k \geq j_1} \left(\mathbb{M}(\infty, \boldsymbol{\eta})^{-1} (I - \Pi_0(\boldsymbol{\eta})) \right)^{k+1-j_1} (\hat{F}_{j_1}(\boldsymbol{\xi}'), 0, \dots, 0)^T, \end{aligned}$$

where the first line corresponds to the component of \mathcal{W}_{j_1} on $\mathbb{E}^s(\infty, \boldsymbol{\eta})$ and the second line to the component on $\mathbb{E}^u(\infty, \boldsymbol{\eta})$. The latter formula defines a solution to (46) and the last thing to prove is that it belongs to ℓ^2 .

The matrix $\mathbb{M}(\infty, \boldsymbol{\eta})$ in (26) depends periodically and analytically on $\boldsymbol{\eta}$. Moreover it has no eigenvalue on \mathbb{S}^1 for any $\boldsymbol{\eta}$ so we have the bounds

$$\|(\mathbb{M}(\infty, \boldsymbol{\eta}) \Pi_0(\boldsymbol{\eta}))^k\| \leq C r^k, \quad \|(\mathbb{M}(\infty, \boldsymbol{\eta})^{-1} (I - \Pi_0(\boldsymbol{\eta})))^k\| \leq C r^k,$$

with $C > 0$ and $r \in (0, 1)$, uniformly with respect to $\boldsymbol{\eta} \in \mathbb{R}^{d-1}$. These bounds express the exponential decay of the stable components of \mathbb{M} and of the inverse of the unstable components. By standard $\ell^1 \star \ell^2$ convolution bounds, we obtain

$$\sum_{j_1 \geq 1} |\mathcal{W}_{j_1}(\boldsymbol{\xi}')|^2 \leq C \sum_{j_1 \geq 1} |\hat{F}_{j_1}(\boldsymbol{\xi}')|^2,$$

with a constant C that is uniform with respect to the frequency ξ' . This means that the inverse Fourier transform of (the appropriate coordinate of) $(\mathcal{W}_{j_1}(\xi'))$ provides with a sequence $w \in \ell^2$ that is a solution to (44). The proof of Proposition 1 is thus complete. \square

In the case $d = 1$, we can also get a solvability result for the reduction of (3) to an interval when one enforces transparent numerical boundary conditions at each end. Namely, we can reproduce the analysis of Section 2 when truncating the initial domain \mathbb{Z} ‘on the right’ rather than ‘on the left’. Introducing some given integer $J \geq p + r + 2$, the solution to (3) with initial data $f^0, \dots, f^s \in \ell^2(\mathbb{Z})$ vanishing outside of the interval $[p + 1, J - r - 1]$ satisfies (11) together with

$$\forall n \in \mathbb{N}, \quad \forall j \geq J, \quad \sum_{m=0}^n \tilde{\Pi}_{n-m} \begin{pmatrix} u_{(j+p, \cdot)}^m \\ \vdots \\ u_{(j+1-r, \cdot)}^m \end{pmatrix} = 0, \quad \text{with} \quad \tilde{\Pi}_n := \delta_{n0} I - \Pi_n.$$

The expression of the matrices $\tilde{\Pi}_n$ comes from the Laurent series expansion of $\Pi^u(z) = I - \Pi^s(s)$. In particular, the restriction of (3) to the interval $[1 - r, J + p]$ reads (recall that we consider the one-dimensional case here):

$$\begin{cases} \sum_{\sigma=0}^{s+1} Q_\sigma u_j^{n+\sigma} = \Delta t F_j^{n+s+1}, & n \geq 0, j = 1, \dots, J, \\ \sum_{m=0}^{n+s+1} \Pi_{n+s+1-m} \begin{pmatrix} u_p^m \\ \vdots \\ u_{1-r}^m \end{pmatrix} = g_\ell^{n+s+1}, & n \geq 0, \\ \sum_{m=0}^{n+s+1} \tilde{\Pi}_{n+s+1-m} \begin{pmatrix} u_{J+p}^m \\ \vdots \\ u_{J+1-r}^m \end{pmatrix} = g_r^{n+s+1}, & n \geq 0, \\ (u_j^0, \dots, u_j^s) = (f_j^0, \dots, f_j^s), & j = 1 - r, \dots, J + p. \end{cases} \quad (47)$$

Up to compatibility conditions for the boundary source terms g_ℓ^n, g_r^n at the left and right ends of the interval, which will take the form (15) (or a similar one at the right end), the main issue for constructing the solution to (47) is to prove existence for the linear problem:

$$\begin{cases} \sum_{\ell=-r}^p a_{\ell, s+1}(\Delta t, \Delta x) u_{j+\ell} = F_j, & j = 1, \dots, J, \\ \Pi_0 \begin{pmatrix} u_p \\ \vdots \\ u_{1-r} \end{pmatrix} = g_\ell, \quad (I - \Pi_0) \begin{pmatrix} u_{J+p} \\ \vdots \\ u_{J+1-r} \end{pmatrix} = g_r. \end{cases}$$

Since the problem is now *finite dimensional*, we are reduced to proving that the only solution to the homogeneous problem with $(F_j)_{j=1, \dots, J} = 0$, $g_\ell = g_r = 0$, is zero. The result is straightforward in the explicit case since the recurrence relation reduces to $u_j = 0$ for $j = 1, \dots, J$ and because of the expression of Π_0 and $I - \Pi_0$, see Theorem 1. We thus turn to the implicit case. With the notation of the proof of Theorem 1, the problem under study reads

$$U_{j+1} = \mathbb{M}(\infty) U_j, \quad j = 1, \dots, J,$$

with the endpoint conditions

$$\Pi_0 U_1 = 0, \quad (I - \Pi_0) U_{J+1} = 0.$$

In the one-dimensional implicit case, Π_0 is the projector on the stable subspace of $\mathbb{M}(\infty)$ parallel to the unstable subspace. Both subspaces form a direct sum of \mathbb{C}^{p+r} and are invariant by $\mathbb{M}(\infty)$. Hence we have $U_1 \in \mathbb{E}^u(\infty)$, $U_{J+1} \in \mathbb{E}^s(\infty)$ and those conditions propagate to any $j = 1, \dots, J+1$ by the recurrence relation $U_{j+1} = \mathbb{M}(\infty) U_j$. The only possibility for all these conditions to hold is to have $U_j = 0$ for all j , which shows unique solvability for (47) and explains why all the problems considered in [BELV16, BMGN16, EA01, ZWH08, ZE06] are indeed solvable. These examples are quickly reviewed in Section 5 below.

4 Strong stability and semigroup estimate. Proof of Theorem 2

Our goal in this Section is to prove Theorem 2. In the first two paragraphs, we show that when the scheme (3) is non-glancing, then (12) is strongly stable and furthermore satisfies the estimate (17). In the last paragraph, we show that the non-glancing condition is necessary for strong stability.

4.1 Strong stability under the non-glancing condition

Let us recall that from now on, we assume that each ratio $\Delta t / \Delta x_i$, $i = 1, \dots, d$ is constant, and that each operator Q_σ depends on the discretization parameters only through these ratios. In addition to Assumptions 1, 2, 4 and 5, we also assume in all this paragraph that the scheme (3) is non-glancing. We can then follow some of the analysis in [Cou15a] and use the following result from [Cou09].

Theorem 3 (Block reduction of \mathbb{M}). *Let Assumptions 1, 2, 4 and 5 be satisfied, and assume that the scheme (3) is non-glancing. Then for all $\underline{z} \in \overline{\mathcal{U}}$ and all $\underline{\eta} \in \mathbb{R}^{d-1}$, there exists an open neighborhood \mathcal{O} of $(\underline{z}, \underline{\eta})$ in $\mathbb{C} \times \mathbb{R}^{d-1}$ and there exists an invertible matrix $T(z, \eta)$ that is holomorphic/analytic with respect to $(z, \eta) \in \mathcal{O}$ such that:*

$$\forall (z, \eta) \in \mathcal{O}, \quad T(z, \eta)^{-1} \mathbb{M}(z, \eta) T(z, \eta) = \begin{pmatrix} \mathbb{M}_1(z, \eta) & & 0 \\ & \ddots & \\ 0 & & \mathbb{M}_L(z, \eta) \end{pmatrix},$$

where the number L of diagonal blocks and the size ν_ℓ of each block \mathbb{M}_ℓ do not depend on $(z, \eta) \in \mathcal{O}$, and where each block satisfies one of the following three properties:

- there exists $\delta > 0$ such that for all $(z, \eta) \in \mathcal{O}$, $\mathbb{M}_\ell(z, \eta)^* \mathbb{M}_\ell(z, \eta) \geq (1 + \delta) I$,
- there exists $\delta > 0$ such that for all $(z, \eta) \in \mathcal{O}$, $\mathbb{M}_\ell(z, \eta)^* \mathbb{M}_\ell(z, \eta) \leq (1 - \delta) I$,
- $\nu_\ell = 1$, \underline{z} and $\mathbb{M}_\ell(\underline{z}, \underline{\eta})$ belong to \mathbb{S}^1 , and $\underline{z} \partial_z \mathbb{M}_\ell(\underline{z}, \underline{\eta}) \overline{\mathbb{M}_\ell(\underline{z}, \underline{\eta})} \in \mathbb{R} \setminus \{0\}$.

A Corollary of Theorem 3 is a (unique) continuation result for the stable and unstable subspaces of $\mathbb{M}(z, \eta)$.

Corollary 1. *Let Assumptions 1, 2, 4 and 5 be satisfied, and assume that the scheme (3) is non-glancing. Then there exists $\delta > 0$ such that the stable and unstable subspaces $\mathbb{E}^s, \mathbb{E}^u$ of \mathbb{M} extend holomorphically/analytically/periodically over $\{\zeta \in \mathbb{C}, |\zeta| > 1 - 2\delta\} \times \mathbb{R}^{d-1}$. In particular, there holds*

$$\sum_{n \geq 0} \frac{1}{(1 - \delta)^n} \sup_{\eta \in \mathbb{R}^{d-1}} \|\Pi_n(\eta)\| < +\infty.$$

Moreover, for all $z \in \mathbb{C}$ with $|z| > 1 - 2\delta$ and $\boldsymbol{\eta} \in \mathbb{R}^{d-1}$, the direct sum (24) holds.

Since the direct sum (24) holds up to $|z| = 1$ (and even a little beyond the unit circle), we can follow the theory in [GKS72], see also [Cou13] for a thorough exposition, and try to prove strong stability for (12) by verifying the so-called Uniform Kreiss-Lopatinskii condition (the main result in [GKS72] is to show that this algebraic condition is actually *equivalent* to strong stability). Verifying the Uniform Kreiss-Lopatinskii condition amounts to first performing a Laplace-Fourier transform in the time and tangential space variables, which reduces (12) to the recurrence relation

$$\begin{cases} \sum_{\ell_1=-r_1}^{p_1} a_{\ell_1}(z, \boldsymbol{\eta}) w_{j_1+\ell_1} = F_{j_1}, & j_1 \geq 1, \\ \Pi^s(z, \boldsymbol{\eta}) \begin{pmatrix} w_{p_1} \\ \vdots \\ w_{1-r_1} \end{pmatrix} = G. \end{cases} \quad (48)$$

Then the goal is to show that when $|z| \geq 1$, there is no non-trivial solution to the homogeneous equation (48) (obtained with $(F_{j_1}) = 0$, $G = 0$). The solutions of interest are those that belong to ℓ^2 when z belongs to \mathcal{U} , and those whose initial data $(w_{p_1}, \dots, w_{1-r_1})^T$ are obtained by a continuation argument from \mathcal{U} to its boundary \mathbb{S}^1 when $z \in \mathbb{S}^1$. For $z \in \mathcal{U}$, we can use Lemma 3 and parametrize the set of ℓ^2 solutions of the recurrence relation

$$\sum_{\ell_1=-r_1}^{p_1} a_{\ell_1}(z, \boldsymbol{\eta}) w_{j_1+\ell_1} = 0, \quad j_1 \geq 1,$$

by the stable subspace $\mathbb{E}^s(z, \boldsymbol{\eta})$ of $\mathbb{M}(z, \boldsymbol{\eta})$. Among all such vectors, it is clear that the only one that satisfies the homogeneous numerical boundary condition

$$\Pi^s(z, \boldsymbol{\eta}) \begin{pmatrix} w_{p_1} \\ \vdots \\ w_{1-r_1} \end{pmatrix} = 0, \quad (49)$$

is the zero vector. In other words, we have just proved that the system (48) has no non-zero solution when the source terms vanish and $z \in \mathcal{U}$. Hence the so-called Godunov-Ryabenkii condition holds for (12) (non-existence of unstable eigenvalues). Proving that the Uniform Kreiss-Lopatinskii condition holds amounts to showing the same ‘injectivity’ property up to $z \in \mathbb{S}^1$. As a first step, let us observe that Corollary 1 shows that the spectral projector Π^s also extends holomorphically/analytically/periodically on a neighborhood of $\overline{\mathcal{U}} \times \mathbb{R}^{d-1}$. It is therefore legitimate to consider the resolvent equation (48) for $z \in \mathbb{S}^1$. In that case, the only vector in the extended stable subspace¹² $\mathbb{E}^s(z, \boldsymbol{\eta})$ that satisfies the homogeneous numerical boundary condition (49) is the zero vector. In other words, we have verified that the Uniform Kreiss-Lopatinskii condition holds. Applying the main result of [GKS72] (more precisely, see [Cou13] for the extension of the theory in [GKS72] to the general case that we consider here), the scheme (12) is strongly stable.

The above argument may look somehow trivial, but the subtle point is that in the theory of [GKS72], one assumes that the numerical boundary conditions for the resolvent equation are ‘well-defined’ for $z \in \overline{\mathcal{U}}$

¹²Recall that for $z \in \mathbb{S}^1$, initial data in $\mathbb{E}^s(z, \boldsymbol{\eta})$ do not necessarily correspond to ℓ^2 solution of the recurrence relation but can be viewed as all the possible limits of such ℓ^2 solutions.

and the difficult part of the job is to extend the stable subspace up to the boundary \mathbb{S}^1 of \mathcal{U} . Here it is not even obvious that the numerical boundary conditions in (48) are well-defined for $z \in \overline{\mathcal{U}}$. As a matter of fact, the main result in [Cou09] shows that the spectral projector Π^s extends (even continuously) up to \mathbb{S}^1 *if and only if* the scheme (3) is non-glancing. When glancing (numerical) wave packets occur, the spectral projector Π^s has a singularity at some point of \mathbb{S}^1 . This singular behavior will be one obstacle we shall have to circumvent in the last paragraph of this Section.

4.2 Semigroup estimate

In this paragraph, we follow the analysis in [CG11, Cou15b] and prove the validity of the stability estimate (17) when one considers non-zero initial data in (12). We therefore assume that the scheme (3) is non-glancing and that for all $\xi \in \mathbb{R}^d$, the roots to (9) are simple (the latter condition being automatic for $s = 0$). Under such assumptions, we can apply the following result from [Cou15b]:

Theorem 4 (Existence of dissipative boundary conditions). *Let Assumptions 1, 2 and 4 be satisfied. Assume furthermore that for all $\xi \in \mathbb{R}^d$, the $s + 1$ roots to (9) are simple. Then there exists a constant $C > 0$ such that for any given initial data $f^0, \dots, f^s \in \ell^2$, there exists a sequence $(v_j^n)_{j_1 \geq 1-r_1, n \in \mathbb{N}}$ that satisfies*

$$\begin{cases} \sum_{\sigma=0}^{s+1} Q_\sigma v_j^{n+\sigma} = 0, & n \geq 0, \quad j_1 \geq 1, \\ (v_j^0, \dots, v_j^s) = (f_j^0, \dots, f_j^s), & j_1 \geq 1 - r_1, \end{cases}$$

and

$$\sup_{n \in \mathbb{N}} \|v^n\|_{1-r_1, +\infty}^2 + \sum_{n \geq s+1} \sum_{j_1=1-r_1}^{p_1} \Delta t \|v_{(j_1, \cdot)}^n\|_{\ell^2(\mathbb{Z}^{d-1})}^2 \leq C \sum_{\sigma=0}^s \|f^\sigma\|_{1-r_1, +\infty}^2.$$

With the help of Theorem 4, we go back to the numerical scheme (12). We tacitly assume again that the source terms in (12) satisfy the necessary compatibility conditions for a solution to exist. We then decompose the solution to (12) as $u_j^n = v_j^n + w_j^n$, where the sequence (v_j^n) is given by Theorem 4, and the remaining part (w_j^n) satisfies

$$\begin{cases} \sum_{\sigma=0}^{s+1} Q_\sigma w_j^{n+\sigma} = \Delta t F_j^{n+s+1}, & n \geq 0, \quad j_1 \geq 1, \\ \sum_{m=0}^{n+s+1} \Pi_{n+s+1-m} \begin{pmatrix} w_{(p_1, \cdot)}^m \\ \vdots \\ w_{(1-r_1, \cdot)}^m \end{pmatrix} = g^{n+s+1} - \sum_{m=0}^{n+s+1} \Pi_{n+s+1-m} \begin{pmatrix} v_{(p_1, \cdot)}^m \\ \vdots \\ v_{(1-r_1, \cdot)}^m \end{pmatrix}, & n \geq 0, \\ (w_j^0, \dots, w_j^s) = (0, \dots, 0), & j_1 \geq 1 - r_1. \end{cases} \quad (50)$$

The main point to keep in mind is that we have reduced to the case of vanishing initial data for (w_j^n) . Since the scheme (3) is non-glancing, we have already seen that (12) is strongly stable and therefore satisfies the estimate (16) (with the interior and boundary source terms as given in (50)). Since we need now to estimate these source terms, we define

$$\forall n \geq 0, \quad \tilde{g}^{n+s+1} := g^{n+s+1} - \sum_{m=0}^{n+s+1} \Pi_{n+s+1-m} \begin{pmatrix} v_{(p_1, \cdot)}^m \\ \vdots \\ v_{(1-r_1, \cdot)}^m \end{pmatrix}.$$

The bound given in Corollary 1 for the matrices $\Pi_n(\eta)$ implies that the Fourier multipliers (Π_n) satisfy

$$\sum_{n \geq 0} \|\Pi_n\|_{\mathcal{B}(\ell^2(\mathbb{Z}^{d-1}))} < +\infty.$$

Actually, the decay is even exponential with respect to n , but we shall only make use of the fact that the norms of these operators belong to ℓ^1 . We use the above definition of the source term \tilde{g}^{n+s+1} and derive the estimates (here $\gamma > 0$ is a parameter and the constants C below are independent of γ):

$$\begin{aligned} \sum_{n \geq s+1} \Delta t e^{-2\gamma n \Delta t} \|\tilde{g}^n\|_{\ell^2(\mathbb{Z}^{d-1})}^2 &\leq C \left\{ \sum_{n \geq s+1} \Delta t e^{-2\gamma n \Delta t} \|g^n\|_{\ell^2(\mathbb{Z}^{d-1})}^2 \right. \\ &\quad \left. + \sum_{n \geq s+1} \Delta t e^{-2\gamma n \Delta t} \left\| \sum_{m=0}^n \Pi_{n-m} \begin{pmatrix} v_{(p_1, \cdot)}^m \\ \vdots \\ v_{(1-r_1, \cdot)}^m \end{pmatrix} \right\|_{\ell^2(\mathbb{Z}^{d-1})}^2 \right\} \\ &\leq C \left\{ \sum_{n \geq s+1} \Delta t e^{-2\gamma n \Delta t} \|g^n\|_{\ell^2(\mathbb{Z}^{d-1})}^2 \right. \\ &\quad \left. + \sum_{n \geq 0} \Delta t e^{-2\gamma n \Delta t} \sum_{j_1=1-r_1}^{p_1} \|v_{(j_1, \cdot)}^n\|_{\ell^2(\mathbb{Z}^{d-1})}^2 \right\} \\ &\leq C \left\{ \sum_{\sigma=0}^s \|f^\sigma\|_{1-r_1, +\infty}^2 + \sum_{n \geq s+1} \Delta t e^{-2\gamma n \Delta t} \|g^n\|_{\ell^2(\mathbb{Z}^{d-1})}^2 \right\}, \end{aligned}$$

where we have first used the standard $\ell^1 \star \ell^2$ convolution estimate and then the bound provided by Theorem 4 for the trace of (v_j^n) .

Since we have an estimate of the source terms in (50), we can use the strong stability of (12) and obtain the following estimate for the solution (w_j^n) to (50):

$$\begin{aligned} &\frac{\gamma}{\gamma \Delta t + 1} \sum_{n \geq 0} \Delta t e^{-2\gamma n \Delta t} \|w^n\|_{1-r_1, +\infty}^2 + \sum_{n \geq s+1} \sum_{j_1=1-r_1}^{p_1} \Delta t e^{-2\gamma n \Delta t} \|w_{(j_1, \cdot)}^n\|_{\ell^2(\mathbb{Z}^{d-1})}^2 \\ &\leq C \left\{ \sum_{\sigma=0}^s \|f^\sigma\|_{1-r_1, +\infty}^2 + \frac{\gamma \Delta t + 1}{\gamma} \sum_{n \geq s+1} \Delta t e^{-2\gamma n \Delta t} \|F^n\|_{1, +\infty}^2 + \sum_{n \geq s+1} \Delta t e^{-2\gamma n \Delta t} \|g^n\|_{\ell^2(\mathbb{Z}^{d-1})}^2 \right\}. \end{aligned} \tag{51}$$

The goal now is to combine (51) with the estimate of Theorem 4 in order to derive the estimate (17) of Theorem 2 we are aiming at. Both (51) and the estimate of Theorem 4 provide with an estimate for the traces of (v_j^n) and (w_j^n) that is sufficient for deriving the estimate of the trace of (u_j^n) in (17). Unfortunately, this is not over yet since on the left hand side of (51), we only control the norm

$$\frac{\gamma}{\gamma \Delta t + 1} \sum_{n \geq 0} \Delta t e^{-2\gamma n \Delta t} \|w^n\|_{1-r_1, +\infty}^2,$$

and not the (stronger) semigroup norm

$$\sup_{n \in \mathbb{N}} e^{-2\gamma n \Delta t} \|w^n\|_{1-r_1, +\infty}^2.$$

However, at this stage, the *exact same* argument as in [Cou15b, Paragraph 3.1] using the multiplier technique developed in that article provides with the semigroup estimate of (w_j^n) . This part of the argument in [Cou15b] is not restricted to the ‘local’ numerical boundary conditions considered in that paper but applies as long as one already controls the trace of the solution to (50) (which is provided here, as in [Cou15b], by the strong stability of (12)). Hence we can improve the estimate (51) into

$$\begin{aligned} & \sup_{n \in \mathbb{N}} e^{-2\gamma n \Delta t} \|w^n\|_{1-r_1, +\infty}^2 + \sum_{n \geq s+1} \sum_{j_1=1-r_1}^{p_1} \Delta t e^{-2\gamma n \Delta t} \|w_{(j_1, \cdot)}^n\|_{\ell^2(\mathbb{Z}^{d-1})}^2 \\ & \leq C \left\{ \sum_{\sigma=0}^s \|f^\sigma\|_{1-r_1, +\infty}^2 + \frac{\gamma \Delta t + 1}{\gamma} \sum_{n \geq s+1} \Delta t e^{-2\gamma n \Delta t} \|F^n\|_{1, +\infty}^2 + \sum_{n \geq s+1} \Delta t e^{-2\gamma n \Delta t} \|g^n\|_{\ell^2(\mathbb{Z}^{d-1})}^2 \right\}, \end{aligned} \quad (52)$$

and combining (52) with the estimate for (v_j^n) provided by Theorem 4, we complete the proof of (17).

4.3 Necessity of the non-glancing condition

Our goal in this Paragraph is to show the last part of Theorem 2, meaning that the non-glancing condition is necessary for strong stability of (12). We therefore assume from now on that the scheme (3) is glancing and that strong stability holds for (12). Our goal will be to obtain a contradiction.

By the analysis of [GKS72], strong stability of (12) is equivalent to the fulfillment of a uniform stability estimate for the solution to the resolvent equation (48). More precisely, since we have assumed that (12) is strongly stable, then there exists a constant $C > 0$ such that for all $z \in \mathcal{U}$ and all $\boldsymbol{\eta} \in \mathbb{R}^{d-1}$, for all $(F_{j_1})_{j_1 \geq 1} \in \ell^2$ and for all $G \in \mathbb{E}^s(z, \boldsymbol{\eta})$, the resolvent equation (48) has a unique solution $(w_{j_1})_{j_1 \geq 1-r_1} \in \ell^2$ and that solution satisfies

$$\frac{|z|-1}{|z|} \sum_{j_1 \geq 1-r_1} |w_{j_1}|^2 + \sum_{j_1=1-r_1}^{p_1} |w_{j_1}|^2 \leq C \left\{ \frac{|z|}{|z|-1} \sum_{j_1 \geq 1} |F_{j_1}|^2 + |G|^2 \right\}.$$

Let us be a little more specific and consider the resolvent equation (48) in the particular case $G = 0$. Then using the companion matrix \mathbb{M} in (22) to rewrite the recurrence relation in (48), we can find an expression for the solution to (48). In particular, there holds

$$\mathcal{W}_1 := \begin{pmatrix} w_{p_1} \\ \vdots \\ w_{1-r_1} \end{pmatrix} = - \sum_{k \geq 1} \mathbb{M}(z, \boldsymbol{\eta})^{-k} \Pi^u(z, \boldsymbol{\eta}) \begin{pmatrix} F_k/a_{p_1}(z, \boldsymbol{\eta}) \\ 0 \\ \vdots \\ 0 \end{pmatrix}. \quad (53)$$

Thanks to our strong stability assumption, we know that uniformly with respect to $(z, \boldsymbol{\eta}) \in \mathcal{U} \times \mathbb{R}^{d-1}$, the vector \mathcal{W}_1 defined in (53) satisfies

$$|\mathcal{W}_1|^2 \leq C \frac{|z|}{|z|-1} \sum_{j_1 \geq 1} |F_{j_1}|^2. \quad (54)$$

The inconvenient feature of (53) is that the source term, meaning the vector to which we apply the matrix $\mathbb{M}(z, \boldsymbol{\eta})^{-k} \Pi^u(z, \boldsymbol{\eta})$ on the right hand side, should be proportional to the first vector of the canonical basis. However, one can argue as in [Cou13, Proposition 4] and show that the same estimate as (54) holds for arbitrary source terms in $\mathbb{C}^{p_1+r_1}$. More precisely, if strong stability holds and under the assumptions of Theorem 2, there exists a constant $C > 0$ such that for all $z \in \mathcal{U}$ with $|z| \leq 2$, for all $\boldsymbol{\eta} \in \mathbb{R}^{d-1}$ and for all $(\mathcal{F}_{j_1})_{j_1 \geq 1} \in \ell^2$, the vector

$$\mathcal{W}_1 := - \sum_{k \geq 1} \mathbb{M}(z, \boldsymbol{\eta})^{-k} \Pi^u(z, \boldsymbol{\eta}) \mathcal{F}_k, \quad (55)$$

satisfies the estimate

$$|\mathcal{W}_1|^2 \leq C \frac{|z|}{|z| - 1} \sum_{j_1 \geq 1} |\mathcal{F}_{j_1}|^2. \quad (56)$$

Our goal now is to show that, if (3) is glancing, then the estimate (56) breaks down for a convenient choice of the frequency z (that should be sufficiently close to \mathcal{U}) and of the source term $(\mathcal{F}_{j_1})_{j_1 \geq 1} \in \ell^2$. Let us therefore recall the following result from [Cou09], which will be the starting point for our construction of the source term $(\mathcal{F}_{j_1})_{j_1 \geq 1} \in \ell^2$ in (55).

Theorem 5 (Block reduction of \mathbb{M}). *Let Assumptions 1, 2, 4 and 5 be satisfied, and assume that the scheme (3) is glancing. Then there exists $\underline{z} \in \overline{\mathcal{U}}$ and $\underline{\boldsymbol{\eta}} \in \mathbb{R}^{d-1}$, there exists an open neighborhood \mathcal{O} of $(\underline{z}, \underline{\boldsymbol{\eta}})$ in $\mathbb{C} \times \mathbb{R}^{d-1}$ and there exists an invertible matrix $T(z, \boldsymbol{\eta})$ that is holomorphic/analytic with respect to $(z, \boldsymbol{\eta}) \in \mathcal{O}$ such that:*

$$\forall (z, \boldsymbol{\eta}) \in \mathcal{O}, \quad T(z, \boldsymbol{\eta})^{-1} \mathbb{M}(z, \boldsymbol{\eta}) T(z, \boldsymbol{\eta}) = \begin{pmatrix} \mathbb{M}_g(z, \boldsymbol{\eta}) & 0 \\ 0 & \mathbb{M}_\#(z, \boldsymbol{\eta}) \end{pmatrix},$$

where the diagonal block \mathbb{M}_g has size $m \times m$, $m \geq 2$, and it satisfies

$$\mathbb{M}_g(\underline{z}, \underline{\boldsymbol{\eta}}) = \underline{\kappa} \begin{pmatrix} 1 & 1 & 0 & 0 \\ 0 & \ddots & \ddots & 0 \\ \vdots & \ddots & \ddots & 1 \\ 0 & \dots & 0 & 1 \end{pmatrix}, \quad \underline{\kappa} \in \mathbb{S}^1.$$

Moreover the lower left coefficient $\underline{\nu}$ of $\partial_z \mathbb{M}_g(\underline{z}, \underline{\boldsymbol{\eta}})$ is such that for all $\theta \in \mathbb{C}$ with $\operatorname{Re} \theta > 0$, and for all complex number ζ such that $\zeta^m = \overline{\underline{\kappa}} \underline{\nu} \underline{z} \theta$, then $\operatorname{Re} \zeta \neq 0$.

In the terminology of [Cou09], the block \mathbb{M}_g is of the fourth type (the subscript g here refers to ‘glancing’). In all what follows, we keep the tangential frequency $\boldsymbol{\eta}$ in (55) fixed and equal to that $\underline{\boldsymbol{\eta}} \in \mathbb{R}^{d-1}$ given in Theorem 5. We shall also choose $z = z_\varepsilon := (1 + \varepsilon) \underline{z}$, with $\varepsilon > 0$ small enough so that $(z_\varepsilon, \underline{\boldsymbol{\eta}})$ belongs to the neighborhood \mathcal{O} of $(\underline{z}, \underline{\boldsymbol{\eta}})$ given by Theorem 5. Since $\boldsymbol{\eta}$ is fixed, we forget to recall the $\boldsymbol{\eta}$ -dependence of all quantities from now on.

The m first column vectors of the matrix $T(z)$ in Theorem 5 are denoted $T_1(z), \dots, T_m(z)$. They satisfy

$$\mathbb{M}(z) \begin{pmatrix} T_1(z) & \dots & T_m(z) \end{pmatrix} = \begin{pmatrix} T_1(z) & \dots & T_m(z) \end{pmatrix} \mathbb{M}_g(z).$$

At $z = \underline{z}$, $\mathbb{M}_g(\underline{z})$ has the (only) eigenvalue $\underline{\kappa} \in \mathbb{S}^1$ with algebraic multiplicity m . The geometric multiplicity is 1. For $z \in \mathcal{U}$ close to \underline{z} , we know from Lemma 3 that $\mathbb{M}_g(z)$ has no eigenvalue on \mathbb{S}^1 for eigenvalues of $\mathbb{M}_g(z)$ are also eigenvalues of $\mathbb{M}(z)$. Therefore the number μ of stable eigenvalues of $\mathbb{M}_g(z)$ when $z \in \mathcal{U}$

is close to \underline{z} is constant. This number is denoted μ from now on. Its value is given in [Cou11, Proposition 4.1]:

$$\mu = \begin{cases} m/2, & \text{if } m \text{ is even,} \\ (m \pm 1)/2, & \text{if } m \text{ is odd.} \end{cases}$$

The choice between ± 1 when m is odd depends on the lower left coefficient $\underline{\nu}$ of $\partial_z \mathbb{M}_g(\underline{z})$. The eigenvalues and eigenvectors of $\mathbb{M}_g(z)$ have a Puiseux expansion close to \underline{z} , see [Bau85]. Such expansions are computed as in [Cou11, Proposition 4.1] (see similar arguments for the continuous problem in [Kre70, Sar65]). The expansions read

$$\begin{aligned} \kappa(z) &= \underline{\kappa} \left(1 + \zeta_1 w^{1/m} + \zeta_2 w^{2/m} + \cdots + \zeta_k w^{k/m} + \cdots \right), \\ r(z) &= \mathbf{r}_0 + \mathbf{r}_1 w^{1/m} + \mathbf{r}_2 w^{2/m} + \cdots + \mathbf{r}_k w^{k/m} + \cdots, \end{aligned}$$

with $w := (z - \underline{z})/\underline{z}$, and the vectors $\mathbf{r}_0, \dots, \mathbf{r}_{m-1}$ form a basis of \mathbb{C}^m . Moreover the first coefficient ζ_1 in the expansion of $\kappa(z)$ is nonzero and can be chosen to be real if m is odd. Independently of m , the number ζ_1 is such that for any m -th root of unity ω , the real part of $\zeta_1 \omega$ is nonzero.

Let us label the m -th roots of unity as $\omega_1, \dots, \omega_m$ and specify $z = z_\varepsilon := (1 + \varepsilon) \underline{z}$, $\varepsilon > 0$ small enough. Then the eigenvalues $\kappa_\ell(\varepsilon)$, $\ell = 1, \dots, m$, of $\mathbb{M}_g(z_\varepsilon)$ have the expansions

$$\kappa_\ell(\varepsilon) = \underline{\kappa} \left(1 + \zeta_1 \omega_\ell \varepsilon^{1/m} \right) + O(\varepsilon^{2/m}),$$

and the associated eigenvectors read

$$r_\ell(\varepsilon) = \sum_{k=0}^{m-1} \mathbf{r}_k \omega_\ell^k \varepsilon^{k/m} + O(\varepsilon).$$

The m -th roots of unity are labeled in such a way that $\kappa_1(\varepsilon), \dots, \kappa_\mu(\varepsilon)$ are the stable eigenvalues of $\mathbb{M}_g(z_\varepsilon)$, and $\kappa_{\mu+1}(\varepsilon), \dots, \kappa_m(\varepsilon)$ are the unstable eigenvalues. To each eigenvector $r_\ell(\varepsilon)$ for $\mathbb{M}_g(z_\varepsilon)$, there corresponds an eigenvector

$$\mathcal{T}_\ell(\varepsilon) := \begin{pmatrix} T_1(z_\varepsilon) & \cdots & T_m(z_\varepsilon) \end{pmatrix} r_\ell(\varepsilon),$$

for the matrix $\mathbb{M}(z_\varepsilon)$, with the same eigenvalue $\kappa_\ell(\varepsilon)$. In particular, $\mathcal{T}_1, \dots, \mathcal{T}_\mu$ are stable eigenvectors and $\mathcal{T}_{\mu+1}, \dots, \mathcal{T}_m$ are unstable eigenvectors.

The goal is now to choose a source term (\mathcal{F}_{j_1}) of size ~ 1 in ℓ^2 but such that the projection on the unstable subspace of $\mathbb{M}(z_\varepsilon)$ is large and proportional to some given unstable eigenvector. Namely, we first define

$$\mathcal{F}(\varepsilon) := \sum_{\ell=1}^{\mu+1} \alpha_\ell(\varepsilon) \mathcal{T}_\ell(\varepsilon) = \begin{pmatrix} T_1(z_\varepsilon) & \cdots & T_m(z_\varepsilon) \end{pmatrix} \sum_{\ell=1}^{\mu+1} \alpha_\ell(\varepsilon) r_\ell(\varepsilon), \quad (57)$$

where the coefficients $\alpha_1(\varepsilon), \dots, \alpha_{\mu+1}(\varepsilon)$ are chosen such that

$$\begin{pmatrix} 1 & \cdots & 1 \\ \vdots & & \vdots \\ \omega_1^\mu & \cdots & \omega_{\mu+1}^\mu \end{pmatrix} \begin{pmatrix} \alpha_1(\varepsilon) \\ \vdots \\ \alpha_{\mu+1}(\varepsilon) \end{pmatrix} = \begin{pmatrix} 0 \\ \vdots \\ 0 \\ \varepsilon^{-\mu/m} \end{pmatrix}. \quad (58)$$

With this definition of the coefficients $\alpha_\ell(\varepsilon)$, and using the above Puiseux expansions of the eigenvectors $r_\ell(\varepsilon)$, we have

$$\mathcal{F}(\varepsilon) = \begin{pmatrix} T_1(\underline{z}) & \cdots & T_m(\underline{z}) \end{pmatrix} \mathbf{r}_\mu + o(1),$$

so, for some positive constant $c > 0$, we have

$$|\mathcal{F}(\varepsilon)| = c + o(1),$$

as ε tends to zero. With $z = z_\varepsilon$, we choose in (55) the source term

$$\forall j_1 \geq 1, \quad \mathcal{F}_{j_1} := \frac{(|\kappa_{\mu+1}(\varepsilon)|^2 - 1)^{1/2}}{\kappa_{\mu+1}(\varepsilon)^{j_1}} \mathcal{F}(\varepsilon),$$

with the vector $\mathcal{F}(\varepsilon)$ given in (57). This choice gives

$$\sum_{j_1 \geq 1} |\mathcal{F}_{j_1}|^2 = |\mathcal{F}(\varepsilon)|^2 = c + o(1),$$

and the corresponding vector \mathcal{W}_1 in (55) reads

$$\begin{aligned} \mathcal{W}_1 &= -\alpha_{\mu+1}(\varepsilon) \sum_{k \geq 1} \frac{(|\kappa_{\mu+1}(\varepsilon)|^2 - 1)^{1/2}}{\kappa_{\mu+1}(\varepsilon)^k} \mathbb{M}(z_\varepsilon)^{-k} \mathcal{T}_{\mu+1}(\varepsilon) \\ &= -\alpha_{\mu+1}(\varepsilon) \sum_{k \geq 1} \frac{(|\kappa_{\mu+1}(\varepsilon)|^2 - 1)^{1/2}}{|\kappa_{\mu+1}(\varepsilon)|^{2k}} \mathcal{T}_{\mu+1}(\varepsilon) \\ &= -\frac{\alpha_{\mu+1}(\varepsilon)}{(|\kappa_{\mu+1}(\varepsilon)|^2 - 1)^{1/2}} \mathcal{T}_{\mu+1}(\varepsilon). \end{aligned}$$

The bound (56) then gives

$$\frac{|\alpha_{\mu+1}(\varepsilon)|^2}{|\kappa_{\mu+1}(\varepsilon)|^2 - 1} |\mathcal{T}_{\mu+1}(\varepsilon)|^2 \leq \frac{C}{\varepsilon},$$

with $C > 0$ uniform with respect to ε . The conclusion follows from the asymptotics of the quantities on the left hand side of this last inequality. Namely, from our construction of the eigenvectors \mathcal{T}_ℓ and the Puiseux expansion of the eigenvalues (recall that the real part of $\zeta_1 \omega_{\mu+1}$ is non-zero and therefore positive since $\kappa_{\mu+1}(\varepsilon)$ is an unstable eigenvalue), we have

$$|\kappa_{\mu+1}(\varepsilon)|^2 - 1 \sim c \varepsilon^{1/m}, \quad |\mathcal{T}_{\mu+1}(\varepsilon)|^2 \sim c,$$

with c a positive constant that does not depend on ε . Eventually, the inverse of the Vandermonde matrix in (58) has a nonzero lower right coefficient (use (30)):

$$\alpha_{\mu+1}(\varepsilon) = \frac{\varepsilon^{-\mu/m}}{\prod_{\ell=1}^{\mu} (\omega_{\mu+1} - \omega_\ell)}.$$

In other words, we have shown that for a suitable constant $C > 0$, and for all $\varepsilon > 0$ arbitrarily small, there holds

$$1 \leq C \varepsilon^{(2\mu+1)/m-1} + o(\varepsilon^{(2\mu+1)/m-1}).$$

Because of the already mentioned result of [Cou11, Proposition 4.1], this forces m to be odd and μ to equal $(m-1)/2$ for otherwise we are led to a contradiction.

It therefore remains to deal with the last possible case: $m (\geq 3)$ is an odd number, and $\mu = (m-1)/2$. In particular, there are at least two unstable eigenvalues for $\mathbb{M}_g(z_\varepsilon)$. The proof of Theorem 2 in this case is a slight refinement of the above argument, which consists in first defining (compare with (57)):

$$\mathcal{F}(\varepsilon) := \sum_{\ell=1}^{\mu+2} \alpha_\ell(\varepsilon) \mathcal{T}_\ell(\varepsilon),$$

with coefficients $\alpha_1(\varepsilon), \dots, \alpha_{\mu+2}(\varepsilon)$ defined by the relation

$$\begin{pmatrix} 1 & \cdots & 1 \\ \vdots & & \vdots \\ \omega_1^\mu & \cdots & \omega_{\mu+2}^\mu \end{pmatrix} \begin{pmatrix} \alpha_1(\varepsilon) \\ \vdots \\ \alpha_{\mu+2}(\varepsilon) \end{pmatrix} = \begin{pmatrix} 0 \\ \vdots \\ 0 \\ \varepsilon^{-(\mu+1)/m} \end{pmatrix}.$$

In that case, we have

$$\Pi^u(z_\varepsilon) \mathcal{F}(\varepsilon) = \alpha_{\mu+1}(\varepsilon) \mathcal{T}_{\mu+1}(\varepsilon) + \alpha_{\mu+2}(\varepsilon) \mathcal{T}_{\mu+2}(\varepsilon),$$

which is a little less nice than before because the unstable part of $\mathcal{F}(\varepsilon)$ has components on two (mostly parallel) eigenvectors of $\mathbb{M}(z_\varepsilon)$ so there might be some cancelation between those two components. We keep nevertheless the same definition as above for the source term in (55), namely:

$$\forall j_1 \geq 1, \quad \mathcal{F}_{j_1} := \left(|\kappa_{\mu+1}(\varepsilon)|^2 - 1 \right)^{1/2} \frac{1}{\kappa_{\mu+1}(\varepsilon)^{j_1}} \mathcal{F}(\varepsilon),$$

but this time with our new definition of $\mathcal{F}(\varepsilon)$. We still have

$$\sum_{j_1 \geq 1} |\mathcal{F}_{j_1}|^2 = c + o(1),$$

but now the corresponding vector \mathcal{W}_1 in (55) is given by

$$- \left(|\kappa_{\mu+1}(\varepsilon)|^2 - 1 \right)^{-1/2} \mathcal{W}_1 = \frac{\alpha_{\mu+1}(\varepsilon)}{|\kappa_{\mu+1}(\varepsilon)|^2 - 1} \mathcal{T}_{\mu+1}(\varepsilon) + \frac{\alpha_{\mu+2}(\varepsilon)}{\kappa_{\mu+1}(\varepsilon) \kappa_{\mu+2}(\varepsilon) - 1} \mathcal{T}_{\mu+2}(\varepsilon).$$

We use again the explicit formula (30) for the inverse of a Vandermonde matrix to derive

$$\begin{aligned} \alpha_{\mu+1}(\varepsilon) &= \frac{\varepsilon^{-(\mu+1)/m}}{(\omega_{\mu+1} - \omega_{\mu+2}) \prod_{k=1}^{\mu} (\omega_{\mu+1} - \omega_k)}, \\ \alpha_{\mu+2}(\varepsilon) &= \frac{\varepsilon^{-(\mu+1)/m}}{(\omega_{\mu+2} - \omega_{\mu+1}) \prod_{k=1}^{\mu} (\omega_{\mu+2} - \omega_k)}, \end{aligned}$$

and the Puiseux expansion of the eigenvalues κ_ℓ give

$$|\kappa_{\mu+1}(\varepsilon)|^2 - 1 \sim \zeta_1 (\overline{\omega_{\mu+1}} + \omega_{\mu+1}) \varepsilon^{1/m}, \quad \overline{\kappa_{\mu+1}(\varepsilon)} \kappa_{\mu+2}(\varepsilon) - 1 \sim \zeta_1 (\overline{\omega_{\mu+1}} + \omega_{\mu+2}) \varepsilon^{1/m}.$$

Moreover, the eigenvectors $\mathcal{T}_{\mu+1}(\varepsilon)$ and $\mathcal{T}_{\mu+2}(\varepsilon)$ share the same finite non zero limit as ε tends to zero. Simplifying the previous expression of \mathcal{W}_1 by the non-zero quantity $\zeta_1 (\omega_{\mu+2} - \omega_{\mu+1})$, we obtain that there exists a non-zero vector \mathbf{W} in $\mathbb{C}^{p_1+r_1}$ such that

$$\varepsilon^{(\mu+3/2)/m} \mathcal{W}_1 \rightarrow \left(\frac{1}{(\overline{\omega_{\mu+1}} + \omega_{\mu+1}) \prod_{k=1}^{\mu} (\omega_{\mu+1} - \omega_k)} - \frac{1}{(\overline{\omega_{\mu+1}} + \omega_{\mu+2}) \prod_{k=1}^{\mu} (\omega_{\mu+2} - \omega_k)} \right) \mathbf{W}.$$

Since $2\mu + 3$ is larger than m , we shall obtain a contradiction as in the previous simpler analysis provided that we can show that the quantity

$$(\overline{\omega_{\mu+1}} + \omega_{\mu+1}) \prod_{k=1}^{\mu} (\omega_{\mu+1} - \omega_k) - (\overline{\omega_{\mu+2}} + \omega_{\mu+2}) \prod_{k=1}^{\mu} (\omega_{\mu+2} - \omega_k), \quad (59)$$

is non-zero.

Let us recall a little the situation here. The ω_ℓ , $\ell = 1, \dots, \mu$ are the m -th roots of unity for which κ_ℓ is a stable eigenvalue. This corresponds to those m -th roots of unity for which $\zeta_1 \omega_\ell$ has negative real part (recall that ζ_1 is a non-zero real number). In particular, the set $\{\omega_1, \dots, \omega_\mu\}$ is invariant by complex conjugation. Furthermore, $\omega_{\mu+1}$ and $\omega_{\mu+2}$ are any m -th roots of unity for which $\zeta_1 \omega_\ell$ has positive real part. In particular, one can always choose $\omega_{\mu+2}$ as the complex conjugate of $\omega_{\mu+1}$. In that case, (59) reduces to showing

$$(\overline{\omega_{\mu+1}} + \omega_{\mu+1}) \prod_{k=1}^{\mu} (\omega_{\mu+1} - \omega_k) - 2\overline{\omega_{\mu+1}} \overline{\prod_{k=1}^{\mu} (\omega_{\mu+1} - \omega_k)} \neq 0,$$

and a sufficient condition for this to happen is

$$(\operatorname{Im} \omega_{\mu+1}) \operatorname{Im} \left(\prod_{k=1}^{\mu} (\omega_{\mu+1} - \omega_k) \right) \neq 0. \quad (60)$$

If ζ_1 is negative, then necessarily m is of the form $3 + 4M$, $M \in \mathbb{N}$, and $\mu = 1 + 2M$. We then choose

$$\omega_{\mu+1} := \exp \left(2i\pi \frac{M+1}{4M+3} \right),$$

and one can rather easily check that the condition (60) is satisfied. A similar calculation yields (60) if ζ_1 is positive (in that case m is of the form $5 + 4M$, $M \in \mathbb{N}$).

5 Examples

5.1 Numerical schemes for the transport equation

In this first Paragraph, the underlying partial differential equation we consider is a one-dimensional transport equation:

$$\partial_t v + a \partial_x v = 0,$$

with $a \neq 0$ a fixed velocity. We explain how the theory developed in this article applies to two possible discretizations of that equation, namely the Lax-Wendroff and leap-frog schemes, both of which being second order in time and space (at least for sufficiently smooth solutions).

The Lax-Wendroff scheme. Given some time and space steps $\Delta t, \Delta x$, and letting for simplicity μ denote the dimensionless parameter

$$\mu := \frac{\Delta t}{\Delta x} a \neq 0,$$

the Lax-Wendroff scheme reads:

$$\begin{cases} u_j^{n+1} - \frac{\mu}{2}(1+\mu)u_{j-1}^n + (\mu^2 - 1)u_j^n + \frac{\mu}{2}(1-\mu)u_{j+1}^n = 0, & n \in \mathbb{N}, j \in \mathbb{Z}, \\ u^0 = f \in \ell^2(\mathbb{Z}). \end{cases} \quad (61)$$

The scheme (61) fits into the framework of (3) with $s = 0$, $Q_1 = I$ (the scheme is explicit), and

$$a_{-1,0} = -\frac{\mu}{2}(1+\mu), \quad a_{0,0} = \mu^2 - 1, \quad a_{1,0} = \frac{\mu}{2}(1-\mu).$$

The stencil of the scheme (61) depends on whether $|\mu| = 1$. To avoid dealing with particular cases, we therefore assume $|\mu| \neq 1$, and therefore $p = r = 1$ in (61). Then it is a rather well-known fact, see e.g. [GKO95], that the Lax-Wendroff scheme satisfies the stability Assumption 2 if and only if $|\mu| < 1$. Since $s = 0$, the stability estimate (6) is even satisfied with the constant 1 (the ℓ^2 norm of the solution to (61) is nonincreasing with respect to n). We thus fix from now on a positive constant λ such that $\lambda|a| < 1$. The set Δ of discretization parameters is

$$\Delta := \{(\Delta t, \Delta t/\lambda), \quad \Delta t \in (0, 1]\},$$

and we consider the iteration (61) with $\mu := \lambda a$. The coefficients in (61) are therefore independent of the discretization parameters $(\Delta t, \Delta x) \in \Delta$. Because $a_{-1,0}$ and $a_{1,0}$ are nonzero, both Assumptions 1 and Assumption 5 are satisfied. We have also seen that Assumption 2 is satisfied thanks to our restriction on λ . The definition (10) reduces here to

$$a_{-1}(z) = a_{-1,0} \neq 0, \quad a_1(z) = a_{1,0} \neq 0,$$

and therefore Assumption 4 is satisfied (the strong version of the non-characteristic boundary assumption). We are now going to verify that the Lax-Wendroff scheme is non-glancing. The amplification matrix \mathcal{A} in (7) reduces here to the complex number

$$\forall \kappa \in \mathbb{C} \setminus \{0\}, \quad \mathcal{A}(\kappa) = \frac{\mu}{2}(1+\mu)\kappa^{-1} - \frac{\mu}{2}(1-\mu)\kappa + 1 - \mu^2.$$

We thus compute

$$\forall \xi \in \mathbb{R}, \quad |\mathcal{A}(e^{i\xi})|^2 = 1 - 4\mu^2(1-\mu^2)\sin^4 \frac{\xi}{2} \leq 1.$$

Therefore the only $\kappa \in \mathbb{S}^1$ for which $\mathcal{A}(\kappa)$ has modulus 1 is $\kappa = 1$, but then $\mathcal{A}'(1) = -\mu$ is nonzero. The Lax-Wendroff scheme (61) is therefore non-glancing¹³ and we can apply Theorem 2 when the computational domain \mathbb{Z} is truncated on one side with the transparent numerical boundary conditions which we now make explicit.

We examine the transparent boundary condition at $j = 0$ associated with (61). For convenience, we adopt here the formulation (37) using a linear form for describing the unstable subspace $\mathbb{E}^u(z)$ of the matrix $\mathbb{M}(z)$ in (36). Let us recall that in (37), the coefficients $\tilde{\kappa}_{u,n}$ correspond to the Laurent series of κ_u^{-1} , where $\kappa_u(z) \in \mathcal{U}$ is the unstable root to the equation

$$a_{-1}(z) + a_0(z)\kappa + a_1(z)\kappa^2 = 0.$$

For the Lax-Wendroff scheme (61), $\kappa_u(z)^{-1}$ satisfies the equation

$$\forall z \in \mathcal{U}, \quad -\mu(1+\mu)\kappa_u(z)^{-2} + 2(z+\mu^2-1)\kappa_u(z)^{-1} + \mu(1-\mu) = 0.$$

¹³This is actually a consequence here of the fact that the scheme is dissipative of order 4, see [GKO95, Chapter 5].

We plug the series $\sum_{n \geq 1} \tilde{\kappa}_{u,n} z^{-n}$ in the latter equation and identify inductively the coefficients, which yields:

$$\begin{aligned} \forall n \geq 2, \quad \tilde{\kappa}_{u,n+1} &= (1 - \mu^2) \tilde{\kappa}_{u,n} + \frac{1}{2} \mu (1 + \mu) \sum_{m=1}^{n-1} \tilde{\kappa}_{u,m} \tilde{\kappa}_{u,n-m}, \\ \text{with } \tilde{\kappa}_{u,1} &= -\frac{1}{2} \mu (1 - \mu), \quad \tilde{\kappa}_{u,2} = (1 - \mu^2) \tilde{\kappa}_{u,1}. \end{aligned} \quad (62)$$

From the relation $\kappa_s(z) \kappa_u(z) = -(1 + \mu)/(1 - \mu)$, we also get the Laurent series expansion of κ_s :

$$\kappa_s(z) = \sum_{n \geq 1} \frac{\kappa_{s,n}}{z^n}, \quad \kappa_{s,n} := -\frac{1 + \mu}{1 - \mu} \tilde{\kappa}_{u,n}. \quad (63)$$

The stable eigenvalue $\kappa_s(z)$ is used to write the transparent numerical boundary condition at the right end $j = J + 1$ of the computational domain.

If one thus truncates the computation domain \mathbb{Z} on both sides and therefore reduces to an interval $[0, J + 1]$, the resulting numerical scheme reads

$$\begin{cases} u_j^{n+1} - \frac{\mu}{2} (1 + \mu) u_{j-1}^n + (\mu^2 - 1) u_j^n + \frac{\mu}{2} (1 - \mu) u_{j+1}^n = 0, & n \in \mathbb{N}, j = 1, \dots, J, \\ u_0^{n+1} = \sum_{m=0}^n \tilde{\kappa}_{u,n+1-m} u_1^m, & n \in \mathbb{N}, \\ u_{J+1}^{n+1} = \sum_{m=0}^n \kappa_{s,n+1-m} u_J^m, & n \in \mathbb{N}, \end{cases}$$

with coefficients $\tilde{\kappa}_{u,n}, \kappa_{s,n}$ defined in (62), (63), and some given initial condition $(u_0^0, \dots, u_{J+1}^0)^T \in \mathbb{R}^{J+2}$. The above numerical scheme is rather easily implemented thanks to the recursive formula for the coefficients $\tilde{\kappa}_{u,n}$ and $\kappa_{s,n}$. But of course there are easier and efficient *local* strategies that are based on absorbing boundary conditions and that may work quite as well, see e.g. [Ehr10, Gol77].

The leap-frog scheme. Given some time and space steps $\Delta t, \Delta x$, and letting again for simplicity μ denote the dimensionless parameter

$$\mu := \frac{\Delta t}{\Delta x} a,$$

the leap-frog scheme reads:

$$\begin{cases} u_j^{n+2} + \mu (u_{j+1}^{n+1} - u_{j-1}^{n+1}) - u_j^n = 0, & n \in \mathbb{N}, j \in \mathbb{Z}, \\ (u^0, u^1) = (f^0, f^1) \in \ell^2(\mathbb{Z})^2. \end{cases} \quad (64)$$

The scheme (64) fits into the framework of (3) with $s = 1$, $Q_2 = -Q_0 = I$ (the scheme is explicit), and

$$a_{-1,1} = -\mu, \quad a_{0,0} = 0, \quad a_{1,0} = \mu, \quad p = r = 1.$$

Assumptions 1 and 5 are thus satisfied. It is also a standard result that the leap-frog scheme is ℓ^2 stable if and only if $|\mu| < 1$, see [RM67]. In that case, Assumption 2 is satisfied with a constant C that only depends on μ . From the above expression of the coefficients, we can also easily check that Assumption 4 is satisfied.

The amplification matrix \mathcal{A} in (7) reads

$$\mathcal{A}(\kappa) = \begin{pmatrix} -\mu(\kappa - \kappa^{-1}) & 1 \\ 1 & 0 \end{pmatrix}.$$

The eigenvalues of $\mathcal{A}(e^{i\xi})$, $\xi \in \mathbb{R}$, are

$$\pm \sqrt{1 - \mu^2 \sin^2 \xi} - i \mu \sin \xi \in \mathbb{S}^1.$$

The derivative of these functions with respect to ξ vanishes when $\xi - \pi/2$ belongs to $\mathbb{Z}\pi$, which means that the leap-frog scheme admits glancing wave packets, which prevents from applying Theorem 2. We can nevertheless derive the transparent boundary conditions by computing the roots $\kappa_s(z) \in \mathbb{D}$ and $\kappa_u(z) \in \mathcal{U}$ to the equation

$$-\mu z \kappa^{-1} + (z^2 - 1) + \mu z \kappa = 0, \quad z \in \mathcal{U}.$$

In particular, the inverse of the unstable root has the following Laurent series expansion

$$\kappa_u(z)^{-1} = \sum_{n \geq 1} \tilde{\kappa}_{u,n} z^{-n},$$

with:

$$\begin{aligned} \forall n \geq 2, \quad \tilde{\kappa}_{u,n+1} &= \tilde{\kappa}_{u,n-1} - \mu \sum_{m=1}^{n-1} \tilde{\kappa}_{u,m} \tilde{\kappa}_{u,n-m}, \\ \tilde{\kappa}_{u,1} &= \mu, \quad \tilde{\kappa}_{u,2} = 0. \end{aligned}$$

The stable root $\kappa_s(z)$ coincides with $-\kappa_u(z)^{-1}$. We can therefore truncate the computation domain \mathbb{Z} and implement the numerical scheme

$$\begin{cases} u_j^{n+2} - \mu(u_{j+1}^{n+1} - u_{j-1}^{n+1}) - u_j^n = 0, & n \in \mathbb{N}, j = 1, \dots, J, \\ u_0^{n+2} = \sum_{m=0}^n \tilde{\kappa}_{u,n+1-m} u_1^m, & n \in \mathbb{N}, \\ u_{J+1}^{n+2} = \sum_{m=0}^n \kappa_{s,n+1-m} u_J^m, & n \in \mathbb{N}, \end{cases}$$

with any given couple of initial data $(u_0^0, \dots, u_{J+1}^0)^T, (u_0^1, \dots, u_{J+1}^1)^T \in \mathbb{R}^{J+2}$.

5.2 Numerical schemes for the heat equation

In this Paragraph, the underlying partial differential equation we consider is the one-dimensional heat equation:

$$\partial_t v - d \partial_{xx}^2 v = 0,$$

with $d > 0$ the diffusion coefficient. We review our derivation of the transparent boundary condition both for the (more than) classical explicit scheme

$$\begin{cases} u_j^{n+1} - \frac{d \Delta t}{\Delta x^2} u_{j-1}^n + 2 \frac{d \Delta t}{\Delta x^2} u_j^n - \frac{d \Delta t}{\Delta x^2} u_{j+1}^n = 0, & n \in \mathbb{N}, j \in \mathbb{Z}, \\ u^0 = f \in \ell^2(\mathbb{Z}), \end{cases} \quad (65)$$

and for the implicit scheme based on the BDF2 quadrature rule (see [HNW93]):

$$\begin{cases} \left(\frac{3}{2} + 2 \frac{d \Delta t}{\Delta x^2} \right) u_j^{n+2} - \frac{d \Delta t}{\Delta x^2} u_{j-1}^{n+2} - \frac{d \Delta t}{\Delta x^2} u_{j+1}^{n+2} - 2 u_j^{n+1} + \frac{1}{2} u_j^n = 0, & n \in \mathbb{N}, j \in \mathbb{Z}, \\ (u^0, u^1) = (f^0, f^1) \in \ell^2(\mathbb{Z})^2. \end{cases} \quad (66)$$

The scheme (65) is of first order in time and second order in space, while (66) is second order in both time and space (again for sufficiently smooth solutions).

Let us first deal with (65), which, in the framework of (3), corresponds to $s = 0$, $Q_1 = I$ (the scheme is explicit) and

$$a_{-1,0}(\Delta t, \Delta x) = a_{1,0}(\Delta t, \Delta x) := -\frac{d \Delta t}{\Delta x^2}, \quad a_{0,0}(\Delta t, \Delta x) := 2 \frac{d \Delta t}{\Delta x^2}.$$

We emphasize here the dependence of the coefficients on Δt and Δx , though of course they only depend on the ratio $\Delta t / \Delta x^2$. The numerical scheme satisfies Assumptions 1, 3 and 5 (it even satisfies the stronger Assumption 4 though we shall not make much use of this fact). As far as the stability Assumption 2 is concerned, it is satisfied if and only if there holds $d \Delta t / \Delta x^2 \leq 1$, which can be readily seen by applying the Fourier transform [RM67] (see also [Str62] for an alternative explanation of this stability condition based on Bernstein's inequality). The admissible set of discretization parameters is therefore

$$\Delta := \{(\Delta t, \Delta x) \in (0, 1]^2 / d \Delta t \leq \Delta x^2\}.$$

For the scheme (65), the derivation of the transparent boundary conditions is based on the analysis of the polynomial equation:

$$a_{-1,0}(\Delta t, \Delta x) + (z + a_{0,0}(\Delta t, \Delta x)) \kappa + a_{1,0}(\Delta t, \Delta x) \kappa^2 = 0,$$

which, in view of the definition of the coefficients $a_{\ell,0}(\Delta t, \Delta x)$, amounts to

$$\mu (\kappa - 1)^2 = z \kappa, \quad \mu := \frac{d \Delta t}{\Delta x^2} \in (0, 1].$$

For $z \in \mathcal{U}$, this equation has one root $\kappa_s(z) \in \mathbb{D}$ and one root $\kappa_u(z) \in \mathcal{U}$, both of which depend holomorphically on z . The Laurent series expansion of $\kappa_u(z)^{-1}$ reads

$$\kappa_u(z)^{-1} = \sum_{n \geq 1} \tilde{\kappa}_{u,n} z^{-n},$$

with:

$$\begin{aligned} \forall n \geq 2, \quad \tilde{\kappa}_{u,n+1} &= -2 \mu \tilde{\kappa}_{u,n} + \mu \sum_{m=1}^{n-1} \tilde{\kappa}_{u,m} \tilde{\kappa}_{u,n-m}, \\ \tilde{\kappa}_{u,1} &= \mu, \quad \tilde{\kappa}_{u,2} = -2 \mu^2. \end{aligned}$$

The stable root $\kappa_s(z)$ coincides with $\kappa_u(z)^{-1}$ so we use the convention $\kappa_{s,n} := \tilde{\kappa}_{u,n}$ in the numerical scheme just below. We can truncate the computation domain \mathbb{Z} in (65) and implement the numerical scheme

$$\begin{cases} u_j^{n+1} - \mu (u_{j-1}^n - 2 u_j^n + u_{j+1}^n) = 0, & n \in \mathbb{N}, j = 1, \dots, J, \\ u_0^{n+1} = \sum_{m=0}^n \tilde{\kappa}_{u,n+1-m} u_1^m, & n \in \mathbb{N}, \\ u_{J+1}^{n+1} = \sum_{m=0}^n \kappa_{s,n+1-m} u_J^m, & n \in \mathbb{N}, \end{cases}$$

with $\mu := d\Delta t/\Delta x^2$ and given initial data $(u_0^0, \dots, u_{J+1}^0)^T \in \mathbb{R}^{J+2}$.

We now consider the implicit scheme (66). Since the scheme is implicit, we first need to determine whether it is well-defined, that is whether Assumption 1 is satisfied. We have

$$\forall \kappa \in \mathbb{C} \setminus \{0\}, \quad \widehat{Q}_2(\kappa) = \frac{3}{2} + 2\mu - \mu(\kappa + \kappa^{-1}),$$

and it is therefore immediate to verify that \widehat{Q}_2 does not vanish on \mathbb{S}^1 . Furthermore, \widehat{Q}_2 has exactly two zeroes that are real (recall $\mu > 0$); one is located in the interval $(0, 1)$ and the remaining one in $(1, +\infty)$. (Actually, one is the inverse of the other.) By the residue Theorem, the index condition (5) is satisfied. Let us now verify Assumption 2. Following [Emm09a, Emm09b], we are going to use the so-called G -stability of the BDF2 quadrature rule, see [HW96, Chapter V.6]. We multiply the recurrence relation in (66) by $\Delta x u_j^{n+2}$ and sum over $j \in \mathbb{Z}$. Using the identity

$$\begin{aligned} 4u_j^{n+2} \left(\frac{3}{2}u_j^{n+2} - 2u_j^{n+1} + \frac{1}{2}u_j^n \right) \\ = (u_j^{n+2})^2 + (2u_j^{n+2} - u_j^{n+1})^2 - (u_j^{n+1})^2 - (2u_j^{n+1} - u_j^n)^2 + (u_j^{n+2} - 2u_j^{n+1} + u_j^n)^2, \end{aligned}$$

and discrete integration by parts in j , we end up with

$$\begin{aligned} \|u^{n+2}\|_{-\infty, +\infty}^2 + \|2u^{n+2} - u^{n+1}\|_{-\infty, +\infty}^2 - \|u^{n+1}\|_{-\infty, +\infty}^2 - \|2u^{n+1} - u^n\|_{-\infty, +\infty}^2 \\ = -\frac{1}{4}\|u^{n+2} - 2u^{n+1} + u^n\|_{-\infty, +\infty}^2 - \frac{\mu}{4}\|(\mathbf{S} - I)u^{n+2}\|_{-\infty, +\infty}^2 \leq 0. \end{aligned}$$

In other words, the energy

$$E^n := \|u^{n+1}\|_{-\infty, +\infty}^2 + \|2u^{n+1} - u^n\|_{-\infty, +\infty}^2,$$

is nonincreasing for solutions to (66), independently of $\mu > 0$, and this proves that Assumption 2 is satisfied with all possible discretization parameters $(\Delta t, \Delta x)$ (that is for $\Delta := (0, 1]^2$, and the corresponding constant C in (6) is independent of $(\Delta t, \Delta x) \in \Delta$). We also compute

$$a_{-1}(z) = -\mu z^2, \quad a_1(z) = -\mu z^2,$$

so Assumptions 3 and 5 are satisfied. We can thus proceed with the construction of transparent boundary conditions for (66). The equation of interest reads

$$\kappa^2 - \frac{1}{\mu z^2} \left\{ \left(\frac{3}{2} + 2\mu \right) z^2 - 2z + \frac{1}{2} \right\} \kappa + 1 = 0.$$

When z belongs to \mathcal{U} , it has one root $\kappa_u(z) \in \mathcal{U}$ and one root $\kappa_s(z) = \kappa_u(z)^{-1} \in \mathbb{D}$. We determine the Laurent series expansion of κ_s :

$$\kappa_s(z) = \sum_{n \geq 0} \frac{\kappa_{s,n}}{z^n}.$$

Plugging this series in the polynomial equation satisfied by κ_s , we end up with the recursive relations:

$$\begin{aligned} \kappa_{s,0}^2 - \frac{1}{\mu} \left(\frac{3}{2} + 2\mu \right) \kappa_{s,0} + 1 &= 0, \quad \kappa_{s,0} \in (0, 1), \\ \left(2\kappa_{s,0} - \frac{3}{2\mu} - 2 \right) \kappa_{s,1} &= -\frac{2\kappa_{s,0}}{\mu}, \\ \forall n \geq 2, \quad \left(2\kappa_{s,0} - \frac{3}{2\mu} - 2 \right) \kappa_{s,n} &= -\frac{2\kappa_{s,n-1}}{\mu} + \frac{\kappa_{s,n-2}}{2\mu} - \sum_{m=1}^{n-1} \kappa_{s,m} \kappa_{s,n-m}. \end{aligned}$$

It should be noted that a crucial fact that we use here is that $\kappa_{s,0}$ is a simple root of the equation $\widehat{Q}_2(\kappa) = 0$, which enables us indeed to determine the sequence $(\kappa_{s,n})$ inductively.

Since the Laurent series expansion of κ_s and κ_u^{-1} coincide, the truncation of (66) on a finite interval $[0, J+1]$ reads

$$\begin{cases} \left(\frac{3}{2} + 2 \frac{d\Delta t}{\Delta x^2} \right) u_j^{n+2} - \frac{d\Delta t}{\Delta x^2} u_{j-1}^{n+2} - \frac{d\Delta t}{\Delta x^2} u_{j+1}^{n+2} - 2u_j^{n+1} + \frac{1}{2} u_j^n = 0, & n \in \mathbb{N}, j = 1, \dots, J, \\ u_0^{n+1} - \kappa_{s,0} u_1^{n+1} = \sum_{m=0}^n \kappa_{s,n+1-m} u_1^m, & n \in \mathbb{N}, \\ u_{J+1}^{n+1} - \kappa_{s,0} u_J^{n+1} = \sum_{m=0}^n \kappa_{s,n+1-m} u_J^m, & n \in \mathbb{N}, \end{cases}$$

with the previous recursive definition for the $\kappa_{s,n}$'s, and any couple of initial conditions $(u_0^0, \dots, u_{J+1}^0)^T, (u_0^1, \dots, u_{J+1}^1)^T \in \mathbb{R}^{J+2}$.

5.3 Numerical schemes for dispersive equations

The two-dimensional Schrödinger equation We consider the two-dimensional linear Schrödinger equation

$$i \partial_t v + \frac{1}{2} \Delta_x v = 0, \quad (t, x) \in \mathbb{R} \times \mathbb{R}^2.$$

We consider the numerical scheme proposed in [EA01] that is based on a centered second order differentiation in space and the Crank-Nicolson quadrature rule. This yields the numerical scheme:

$$\begin{aligned} i \frac{u_{j_1, j_2}^{n+1} - u_{j_1, j_2}^n}{\Delta t} + \frac{1}{4 \Delta x_1^2} \left(u_{j_1+1, j_2}^{n+1} - 2u_{j_1, j_2}^{n+1} + u_{j_1-1, j_2}^{n+1} + u_{j_1+1, j_2}^n - 2u_{j_1, j_2}^n + u_{j_1-1, j_2}^n \right) \\ + \frac{1}{4 \Delta x_2^2} \left(u_{j_1, j_2+1}^{n+1} - 2u_{j_1, j_2}^{n+1} + u_{j_1, j_2-1}^{n+1} + u_{j_1, j_2+1}^n - 2u_{j_1, j_2}^n + u_{j_1, j_2-1}^n \right) = 0, \end{aligned} \quad (67)$$

with some given initial condition $u^0 \in \ell^2(\mathbb{Z}^2; \mathbb{C})$. For future use, we introduce the positive parameters:

$$\mu_1 := \frac{\Delta t}{4 \Delta x_1^2}, \quad \mu_2 := \frac{\Delta t}{4 \Delta x_2^2}.$$

The scheme (67) fits into the framework of (3) with $p_1 = p_2 = r_1 = r_2 = 1$, and with the operators

$$\begin{aligned} Q_1 &:= i I + \mu_1 (\mathbf{S}_1 + \mathbf{S}_1^{-1} - 2 I) + \mu_2 (\mathbf{S}_2 + \mathbf{S}_2^{-1} - 2 I), \\ Q_0 &:= -i I + \mu_1 (\mathbf{S}_1 + \mathbf{S}_1^{-1} - 2 I) + \mu_2 (\mathbf{S}_2 + \mathbf{S}_2^{-1} - 2 I). \end{aligned}$$

In particular, there holds

$$\widehat{Q}_1(e^{i\eta_1}, e^{i\eta_2}) = i - 4\mu_1 \sin^2 \frac{\eta_1}{2} - 4\mu_2 \sin^2 \frac{\eta_2}{2},$$

so that not only $\widehat{Q}_1(e^{i\eta_1}, e^{i\eta_2})$ is nonzero (that is Q_1 is an isomorphism on ℓ^2), but $\widehat{Q}_1(\cdot, e^{i\eta_2})$ maps the unit circle \mathbb{S}^1 into the upper half-plane $\{\zeta \in \mathbb{C}, \operatorname{Im} \zeta > 0\}$. Hence we can write

$$\frac{1}{2i\pi} \int_{\mathbb{S}^1} \frac{\partial_{\kappa_1} \widehat{Q}_1(\kappa_1, e^{i\eta_2})}{\widehat{Q}_1(\kappa_1, e^{i\eta_2})} d\kappa_1 = \frac{1}{2i\pi} \int_{\mathbb{S}^1} \partial_{\kappa_1} (\ln \widehat{Q}_1(\kappa_1, e^{i\eta_2})) d\kappa_1 = 0,$$

where we have used the principal determination of the logarithm. Hence Assumption 1 is satisfied. It is a rather standard property that (67) preserves the ℓ^2 norm, so Assumption 2 is satisfied with the maximal set of discretization parameters $\Delta := (0, 1]^3$ (the constant C in (6) can be chosen to be 1). With the definition (10), we compute (here $\eta = \eta_2$ belongs to \mathbb{R} so we rather use the more explicit notation η_2):

$$a_{-1}(z, \eta_2) = a_1(z, \eta_2) = (z + 1) \mu_1,$$

so Assumption 3 is satisfied (but Assumption 4 is not !). We can also easily check that Assumption 5 is satisfied since a_{-1} and a_1 have degree 1 in z for all η_2 . The derivation of transparent numerical boundary conditions for (67) was performed in [EA01] so we shall not reproduce it here. Approximate (namely, absorbing) numerical boundary conditions for (67) are proposed and studied in [EA01, AES03, AAB⁺08, DZ06]. We also refer to [Sze04, AAB⁺08] and references therein for the construction of absorbing boundary conditions for the nonlinear Schrödinger equation.

The Airy equation We are now going back to a one-dimensional problem and consider as in [ZWH08] the Airy equation

$$\partial_t v + \partial_{xx}^3 v = 0.$$

The extension to a nonzero first order transport term is considered in [BELV16] in view of later dealing with the Korteweg - de Vries equation. For simplicity, we restrict to this simple framework ($U_1 = 0$ and $U_2 = 1$ in the notation of [BELV16]) and explain how one of the schemes considered in [BELV16] fits into our framework. More precisely, we consider the so-called ‘rightside Crank-Nicolson’ scheme proposed in [Qin83]:

$$\frac{u_j^{n+1} - u_j^n}{\Delta t} + \frac{1}{2\Delta x^3} \left(u_{j+2}^{n+1} - 3u_{j+1}^{n+1} + 3u_j^{n+1} - u_{j-1}^{n+1} + u_{j+2}^n - 3u_{j+1}^n + 3u_j^n - u_{j-1}^n \right) = 0, \quad (68)$$

with some given initial condition $u^0 \in \ell^2$. The other scheme considered in [BELV16] is centered in space so the stability analysis is identical to the above one for the discretized Schrödinger equation. With the notation

$$\mu := \frac{\Delta t}{2\Delta x^3} > 0,$$

the numerical scheme (68) fits into the framework of (3) with $p = 2$, $r = 1$, and

$$Q_1 := I + \mu (\mathbf{S}^2 - 3\mathbf{S}_1 + 3I - \mathbf{S}^{-1}), \quad Q_0 := -I + \mu (\mathbf{S}^2 - 3\mathbf{S}_1 + 3I - \mathbf{S}^{-1}).$$

In particular, there holds

$$\widehat{Q_1}(e^{i\xi}) = 1 + 8\mu \sin^4 \frac{\xi}{2} - 4i\mu \sin \xi \sin^2 \frac{\xi}{2}.$$

Hence not only $\widehat{Q_1}(e^{i\xi})$ is nonzero but its real part is not smaller than 1. In particular, we can apply the same argument as above for the discretized Schrödinger equation and use the principal determination of the logarithm to show that $\widehat{Q_1}$ satisfies the index condition (5). It is also proved in [Qin83] that the scheme (68) is stable, that is, it satisfies Assumption 2 with all possible discretization parameters, that is with $\Delta := (0, 1]^2$. Since (68) is based on the Crank-Nicolson quadrature rule, there is no difficulty in verifying that Assumptions 3 and 5 are satisfied. Again Assumption 4 is not satisfied. We refer to [BELV16] for a derivation of the transparent boundary conditions for the scheme (68). In particular, we recover here in our general framework the separation property for the roots $\kappa(z)$ (Theorem 3.1 in [BELV16]). The analysis in [BELV16] makes clear how, for this case with $p = 2$, one can write transparent boundary

conditions by using linear forms rather than projectors. This requires in [BELV16] using some suitable combinations of the two unstable roots (named $\ell_2(z)$ and $\ell_3(z)$ here) in order to preserve holomorphy with respect to z on \mathcal{U} .

The linearized Benjamin-Bona-Mahony equation We discuss eventually the linearized Benjamin-Bona-Mahony equation considered in the recent work [BMGN16]. The partial differential equation under consideration reads

$$\partial_t (u - \varepsilon \partial_{xx}^2 u) + c \partial_x u = 0,$$

with $\varepsilon > 0$ and $c > 0$ (the sign of c is crucial in the ‘upwinding’ procedure below). One numerical scheme proposed in [BMGN16] is based on a centered difference for the dispersive part and on an upwind procedure for the transport part. Then one applies the Crank-Nicolson quadrature rule to integrate in time. The resulting scheme reads:

$$u_j^{n+1} - u_j^n - \frac{\varepsilon \Delta t}{2 \Delta x^2} (u_{j+1}^{n+1} - 2u_j^{n+1} + u_{j-1}^{n+1} - u_{j+1}^n + 2u_j^n - u_{j-1}^n) + \frac{c \Delta t}{2 \Delta x} (u_j^{n+1} - u_{j-1}^{n+1} + u_j^n - u_{j-1}^n) = 0, \quad (69)$$

which fits into the framework of (3) with $p = r = 1$, and

$$Q_1 = - \left(\frac{\varepsilon \Delta t}{2 \Delta x^2} + \frac{c \Delta t}{2 \Delta x} \right) \mathbf{S}^{-1} + \left(1 + \frac{\varepsilon \Delta t}{\Delta x^2} + \frac{c \Delta t}{2 \Delta x} \right) I - \frac{\varepsilon \Delta t}{2 \Delta x^2} \mathbf{S}.$$

In particular, one computes

$$\operatorname{Re} \widehat{Q_1}(e^{i\eta}) = 1 + 2 \frac{\varepsilon \Delta t}{\Delta x^2} \sin^2 \frac{\eta}{2} + \frac{c \Delta t}{\Delta x} \sin^2 \frac{\eta}{2} \geq 1,$$

so not only Q_1 is an isomorphism on $\ell^2(\mathbb{Z})$ but the index condition (5) is satisfied (we use the same argument as for the discretization of the Airy equation). Reproducing more or less the same computations as for (68), one can also show that Assumption 2 is satisfied with all possible discretization parameters, that is with $\Delta = (0, 1]^2$. From the above expression of Q_1 , one can also easily verify that Assumption 5 is satisfied.

We now compute

$$a_{-1}(z) = -(z-1) \frac{\varepsilon \Delta t}{2 \Delta x^2} - (z+1) \frac{c \Delta t}{2 \Delta x}, \quad a_1(z) = -(z-1) \frac{\varepsilon \Delta t}{2 \Delta x^2},$$

so obviously a_1 does not vanish on \mathcal{U} (but it vanishes on $\overline{\mathcal{U}}$). The only root of a_{-1} belongs to \mathbb{D} so Assumption 3 is satisfied (but not Assumption 4). Hence the calculations made explicit in [BMGN16] also fit in the general framework that we have discussed in this article.

Acknowledgments The author warmly thanks Christophe Besse, Vincent Colin, Paolo Ghiggini, François Laudenbach, Laurent Meersseman and Pascal Noble for very helpful and stimulating discussions related to this work.

References

- [AAB⁺08] X. Antoine, A. Arnold, C. Besse, M. Ehrhardt, and A. Schädle. A review of transparent and artificial boundary conditions techniques for linear and nonlinear Schrödinger equations. *Commun. Comput. Phys.*, 4(4):729–796, 2008.

- [AB01] X. Antoine and C. Besse. Construction, structure and asymptotic approximations of a microdifferential transparent boundary conditions for the linear Schrödinger equation. *J. Math. Pures Appl.*, 80(7):701–738, 2001.
- [ABS09] X. Antoine, C. Besse, and J. Szeftel. Towards accurate artificial boundary conditions for nonlinear PDEs through examples. *Cubo*, 11(4):29–48, 2009.
- [AES03] A. Arnold, M. Ehrhardt, and I. Sofronov. Discrete transparent boundary conditions for the Schrödinger equation: fast calculation, approximation, and stability. *Commun. Math. Sci.*, 1(3):501–556, 2003.
- [Aud12] C. Audiard. Non-homogeneous boundary value problems for linear dispersive equations. *Comm. Partial Differential Equations*, 37(1):1–37, 2012.
- [Bau85] H. Baumgärtel. *Analytic perturbation theory for matrices and operators*. Birkhäuser Verlag, 1985.
- [BELV16] C. Besse, M. Ehrhardt, and I. Lacroix-Violet. Discrete artificial boundary conditions for the Korteweg-de-Vries equation. *Numer. Methods Partial Differential Equations*, 2016.
- [BGS07] S. Benzoni-Gavage and D. Serre. *Multidimensional hyperbolic partial differential equations*. Oxford University Press, 2007. First-order systems and applications.
- [BMGN16] C. Besse, B. Mésognon-Gireau, and P. Noble. Artificial boundary conditions for the linearized Benjamin-Bona-Mahony equation. Available at <https://hal.archives-ouvertes.fr/hal-01305360>, 2016.
- [CG11] J.-F. Coulombel and A. Gloria. Semigroup stability of finite difference schemes for multidimensional hyperbolic initial boundary value problems. *Math. Comp.*, 80(273):165–203, 2011.
- [Cou09] J.-F. Coulombel. Stability of finite difference schemes for hyperbolic initial boundary value problems. *SIAM J. Numer. Anal.*, 47(4):2844–2871, 2009.
- [Cou11] J.-F. Coulombel. Stability of finite difference schemes for hyperbolic initial boundary value problems II. *Ann. Sc. Norm. Super. Pisa Cl. Sci. (5)*, X(1):37–98, 2011.
- [Cou13] J.-F. Coulombel. Stability of finite difference schemes for hyperbolic initial boundary value problems. In *HCDTE Lecture Notes. Part I. Nonlinear Hyperbolic PDEs, Dispersive and Transport Equations*, pages 97–225. American Institute of Mathematical Sciences, 2013.
- [Cou15a] J.-F. Coulombel. Fully discrete hyperbolic initial boundary value problems with nonzero initial data. *Confluentes Math.*, 7(2):17–47, 2015.
- [Cou15b] J.-F. Coulombel. The Leray-Gårding method for finite difference schemes. *J. Éc. polytech. Math.*, 2:297–331, 2015.
- [DZ06] B. Ducomet and A. Zlotnik. On stability of the Crank-Nicolson scheme with approximate transparent boundary conditions for the Schrödinger equation. I. *Commun. Math. Sci.*, 4(4):741–766, 2006.

- [EA01] M. Ehrhardt and A. Arnold. Discrete transparent boundary conditions for the Schrödinger equation. *Riv. Mat. Univ. Parma (6)*, 4*:57–108, 2001. Fluid dynamic processes with inelastic interactions at the molecular scale (Torino, 2000).
- [Ehr10] M. Ehrhardt. Absorbing boundary conditions for hyperbolic systems. *Numer. Math. Theory Methods Appl.*, 3(3):295–337, 2010.
- [Emm09a] E. Emmrich. Convergence of the variable two-step BDF time discretisation of nonlinear evolution problems governed by a monotone potential operator. *BIT*, 49(2):297–323, 2009.
- [Emm09b] E. Emmrich. Two-step BDF time discretisation of nonlinear evolution problems governed by monotone operators with strongly continuous perturbations. *Comput. Methods Appl. Math.*, 9(1):37–62, 2009.
- [GF74] I. C. Gohberg and I. A. Fel’dman. *Convolution equations and projection methods for their solution*. American Mathematical Society, 1974. Translated from the Russian, Translations of Mathematical Monographs, Vol. 41.
- [GKO95] B. Gustafsson, H.-O. Kreiss, and J. Oliger. *Time dependent problems and difference methods*. John Wiley & Sons, 1995.
- [GKS72] B. Gustafsson, H.-O. Kreiss, and A. Sundström. Stability theory of difference approximations for mixed initial boundary value problems. II. *Math. Comp.*, 26(119):649–686, 1972.
- [Gol77] M. Goldberg. On a boundary extrapolation theorem by Kreiss. *Math. Comp.*, 31(138):469–477, 1977.
- [GT81] M. Goldberg and E. Tadmor. Scheme-independent stability criteria for difference approximations of hyperbolic initial-boundary value problems. II. *Math. Comp.*, 36(154):603–626, 1981.
- [Hag99] T. Hagstrom. Radiation boundary conditions for the numerical simulation of waves. In *Acta numerica, 1999*, volume 8 of *Acta Numer.*, pages 47–106. Cambridge Univ. Press, 1999.
- [Hal82] L. Halpern. Absorbing boundary conditions for the discretization schemes of the one-dimensional wave equation. *Math. Comp.*, 38(158):415–429, 1982.
- [HNW93] E. Hairer, S. P. Nørsett, and G. Wanner. *Solving ordinary differential equations. I*. Springer-Verlag, second edition, 1993. Nonstiff problems.
- [HW96] E. Hairer and G. Wanner. *Solving ordinary differential equations. II*. Springer-Verlag, second edition, 1996. Stiff and differential-algebraic problems.
- [HY07] H. Han and D. Yin. Absorbing boundary conditions for the multidimensional Klein-Gordon equation. *Commun. Math. Sci.*, 5(3):743–764, 2007.
- [Kat95] T. Kato. *Perturbation theory for linear operators*. Classics in Mathematics. Springer-Verlag, 1995.
- [Kre68] H.-O. Kreiss. Stability theory for difference approximations of mixed initial boundary value problems. I. *Math. Comp.*, 22:703–714, 1968.

- [Kre70] H.-O. Kreiss. Initial boundary value problems for hyperbolic systems. *Comm. Pure Appl. Math.*, 23:277–298, 1970.
- [Lax02] P. D. Lax. *Functional analysis*. Pure and Applied Mathematics. Wiley-Interscience, 2002.
- [Nik02] N. K. Nikolski. *Operators, functions, and systems: an easy reading. Vol. 1*. Mathematical Surveys and Monographs. American Mathematical Society, 2002.
- [Osh69] S. Osher. Systems of difference equations with general homogeneous boundary conditions. *Trans. Amer. Math. Soc.*, 137:177–201, 1969.
- [Osh72] S. Osher. Stability of parabolic difference approximations to certain mixed initial boundary value problems. *Math. Comp.*, 26:13–39, 1972.
- [Qin83] M. Z. Qin. Difference schemes for the dispersive equation. *Computing*, 31(3):261–267, 1983.
- [RM67] R. D. Richtmyer and K. W. Morton. *Difference methods for initial value problems*. Graduate Texts in Mathematics. Interscience Publishers John Wiley & Sons, 1967. Theory and applications.
- [Rud87] W. Rudin. *Real and complex analysis*. McGraw-Hill, 1987.
- [Sar65] L. Sarason. On hyperbolic mixed problems. *Arch. Rational Mech. Anal.*, 18:310–334, 1965.
- [Str62] G. Strang. Trigonometric polynomials and difference methods of maximum accuracy. *J. Math. Phys.*, 41:147–154, 1962.
- [Str64] G. Strang. Wiener-Hopf difference equations. *J. Math. Mech.*, 13:85–96, 1964.
- [SW97] J. C. Strikwerda and B. A. Wade. A survey of the Kreiss matrix theorem for power bounded families of matrices and its extensions. In *Linear operators (Warsaw, 1994)*, volume 38 of *Banach Center Publ.*, pages 339–360. Polish Acad. Sci., 1997.
- [Sze04] J. Szeftel. Design of absorbing boundary conditions for Schrödinger equations in \mathbb{R}^d . *SIAM J. Numer. Anal.*, 42(4):1527–1551, 2004.
- [Sze06] J. Szeftel. Absorbing boundary conditions for the nonlinear Schrödinger equation. *Numer. Math.*, 103:103–127, 2006.
- [Tre84] L. N. Trefethen. Instability of difference models for hyperbolic initial boundary value problems. *Comm. Pure Appl. Math.*, 37:329–367, 1984.
- [VB82] R. Vichnevetsky and J. B. Bowles. *Fourier analysis of numerical approximations of hyperbolic equations*, volume 5 of *SIAM Studies in Applied Mathematics*. Society for Industrial and Applied Mathematics, 1982.
- [ZE06] A. Zisowsky and M. Ehrhardt. Discrete transparent boundary conditions for parabolic systems. *Math. Comput. Modelling*, 43(3-4):294–309, 2006.
- [ZWH08] C. Zheng, X. Wen, and H. Han. Numerical solution to a linearized KdV equation on unbounded domain. *Numer. Methods Partial Differential Equations*, 24(2):383–399, 2008.